

Lab 2

Jacob Minkin

11:59PM February 25, 2021

More Basic R Skills

- 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! (Yes, there is a base R function that does this called `rev`). Use `head` on `v` and `tail` on `my_reverse(v)` to verify it works.

```
my_reverse = function(v) {  
  v_rev = rep(NA, times = length(v))  
  
  for (i in length(v):1) v_rev[length(v)-i+1] = v[i]  
  
  v_rev  
}  
v = 1:10  
my_reverse(v)
```

```
## [1] 10 9 8 7 6 5 4 3 2 1
```

- Create a function `flip_matrix` which takes as required input a matrix, an argument `dim_to_rev` that returns the matrix with the rows in reverse order or the columns in reverse order depending on the `dim_to_rev` argument. Let the default be the dimension of the matrix that is greater.

```
flip_matrix = function(x, dim_to_rev = NULL) {  
  
  if(is.null(dim_to_rev)){  
    dim_to_rev = ifelse(nrow(x) >= ncol(x), "rows", "cols")  
  }  
  
  if(dim_to_rev == "rows"){  
    x[my_reverse(1:nrow(x)),]  
  }  
  else if (dim_to_rev == "cols"){  
    x[,my_reverse(1:ncol(x))]  
  }  
  else stop ("Illegal arg")  
  
}  
  
x = matrix(rnorm(100), nrow=25)  
flip_matrix(x, dim_to_rev = "cols")
```

```
##           [,1]      [,2]      [,3]      [,4]
## [1,]  0.07992786 -0.55007340 -0.77094820  1.3122645686
## [2,] -0.52309123  0.11468824  0.29776821  1.2151784447
## [3,] -1.16377687 -0.62703663  1.04149513 -1.1758543451
## [4,] -0.29191794  1.08034035 -0.17938295 -0.3636391913
## [5,]  0.27111599 -0.91597176  1.70972199 -0.0335488314
## [6,]  1.39287759  0.09416875  0.23054751 -0.0815740850
## [7,] -0.65899518  0.10319489  1.15891071  1.0826073997
## [8,] -0.66051265  0.58762562 -0.08541064 -0.9306360751
## [9,]  0.27376153  0.94573136  0.53997426 -0.5521422849
## [10,] 1.49331772 -2.97830559 -1.12611753  1.0586337492
## [11,] -1.56284958  1.73568261  0.39846352 -0.2594454892
## [12,]  1.46091153  1.04748990 -0.82288546  2.0913818817
## [13,] -0.48518833 -0.72362927 -0.17157304 -0.5521214642
## [14,]  0.08465194 -1.47590952  0.62435894 -0.5913783577
## [15,]  0.01543354  1.25991404  0.06424176  0.6967225341
## [16,] -0.41627276  0.42884839 -0.10323268  1.1210997169
## [17,]  1.29877780 -1.19802789  0.91421060 -1.3459986295
## [18,]  0.05232062 -0.87876155 -0.56772371  0.6853325496
## [19,] -0.21040121 -0.42214842 -0.50783382  0.0001429829
## [20,]  0.51750275 -1.52271389  1.79891048  0.1025119283
## [21,]  0.17951292 -0.35166756  1.62767738 -1.0964398476
## [22,] -0.09132716  0.57082616 -0.31712276  0.3936643098
## [23,]  0.39377100  0.74115270  0.23511357 -1.1464520831
## [24,] -0.01349294  1.57154257 -1.44795382 -1.9665681077
## [25,]  0.62391372  0.56834210 -0.55510958 -0.8350350468
```

- 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 according to this sequence.

```
my_list = list()
for (i in 1:8)
  my_list[LETTERS[i]] = array(1:i^i,dim = c(rep(i,times = i)))
```

```
## Warning in my_list[LETTERS[i]] <- array(1:i^i, dim = c(rep(i, times = i))):
## number of items to replace is not a multiple of replacement length

## Warning in my_list[LETTERS[i]] <- array(1:i^i, dim = c(rep(i, times = i))):
## number of items to replace is not a multiple of replacement length

## Warning in my_list[LETTERS[i]] <- array(1:i^i, dim = c(rep(i, times = i))):
## number of items to replace is not a multiple of replacement length

## Warning in my_list[LETTERS[i]] <- array(1:i^i, dim = c(rep(i, times = i))):
## number of items to replace is not a multiple of replacement length

## Warning in my_list[LETTERS[i]] <- array(1:i^i, dim = c(rep(i, times = i))):
## number of items to replace is not a multiple of replacement length

## Warning in my_list[LETTERS[i]] <- array(1:i^i, dim = c(rep(i, times = i))):
## number of items to replace is not a multiple of replacement length
```

Run the following code:

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

```
## $A
## 56 bytes
##
## $B
## 56 bytes
##
## $C
## 56 bytes
##
## $D
## 56 bytes
##
## $E
## 56 bytes
##
## $F
## 56 bytes
##
## $G
## 56 bytes
##
## $H
## 56 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?

`object.size` provides an estimate of the memory being used to store the object. In the line above it breaks down the list into its different sections so you don't need to input the dimensions of the array.

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

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

A little about strings

- Use the `strsplit` function and `sample` to put the sentences in the string `lorem` below in random order. You will also need to manipulate the output of `strsplit` which is a list. You may need to learn basic concepts of regular expressions.

```
lorem = "Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi posuere varius volutpat. Morbi faucibus ligula id massa ultricies viverra. Donec vehicula sagittis nisi non semper. Donec at tempor erat. Integer dapibus mi lectus, eu posuere arcu ultricies in. Cras suscipit id nibh lacinia elementum. Curabitur est augue, congue eget quam in, scelerisque semper magna. Aenean nulla ante, iaculis sed vehicula ac, finibus vel arcu. Mauris at sodales augue. "
```

```
sample(unlist(strsplit(lorem, "[.]")))
```

```
## [1] " "
## [2] " Mauris at sodales augue"
## [3] " Curabitur est augue, congue eget quam in, scelerisque semper magna"
## [4] " Morbi faucibus ligula id massa ultricies viverra"
## [5] " Aenean nulla ante, iaculis sed vehicula ac, finibus vel arcu"
## [6] " Integer dapibus mi lectus, eu posuere arcu ultricies in"
## [7] "Lorem ipsum dolor sit amet, consectetur adipiscing elit"
## [8] " Morbi posuere varius volutpat"
## [9] " Cras suscipit id nibh lacinia elementum"
## [10] " Donec at tempor erat"
## [11] " Donec vehicula sagittis nisi non semper"
```

You have a set of names divided by gender (M / F) and generation (Boomer / GenX / Millennial):

- M / Boomer "Theodore, Bernard, Gene, Herbert, Ray, Tom, Lee, Alfred, Leroy, Eddie"
- M / GenX "Marc, Jamie, Greg, Darryl, Tim, Dean, Jon, Chris, Troy, Jeff"
- M / Millennial "Zachary, Dylan, Christian, Wesley, Seth, Austin, Gabriel, Evan, Casey, Luis"
- F / Boomer "Gloria, Joan, Dorothy, Shirley, Betty, Dianne, Kay, Marjorie, Lorraine, Mildred"
- F / GenX "Tracy, Dawn, Tina, Tammy, Melinda, Tamara, Tracey, Colleen, Sherri, Heidi"
- F / Millennial "Samantha, Alexis, Brittany, Lauren, Taylor, Bethany, Latoya, Candice, Brittney, Cheyenne"

Create a list-within-a-list that will intelligently store this data.

```

Generations_Gender_list = list(Male =
                                list(Boomer = strsplit("Theodore, Bernard, Gene, Herbert, Ray, Tom, Lee, A
                                lfred, Leroy, Eddie", split = ", ")[[1]],
                                Genx = strsplit("Marc, Jamie, Greg, Darryl, Tim, Dean, Jon, Chris, Tr
                                oy, Jeff", split = ", ")[[1]],
                                Millennial = strsplit("Zachary, Dylan, Christian, Wesley, Seth, Austi
                                n, Gabriel, Evan, Casey, Luis", split = ", ")[[1]]
                                ),
                                Female =
                                list(Boomer = strsplit("Gloria, Joan, Dorothy, Shirley, Betty, Dianne, Kay
                                , Marjorie, Lorraine, Mildred", split = ", ")[[1]],
                                Genx = strsplit("Tracy, Dawn, Tina, Tammy, Melinda, Tamara, Tracey, C
                                olleen, Sherri, Heidi", split = ", ")[[1]],
                                Millennial = strsplit("Samantha, Alexis, Brittany, Lauren, Taylor, Be
                                thany, Latoya, Candice, Brittney, Cheyenne", split = ", ")[[1]]
                                )
                                )
#strsplit("Theodore, Bernard, Gene, Herbert, Ray, Tom, Lee, Alfred, Leroy, Eddie", split = ", ")[[1]]
cat (str(Generations_Gender_list))

```

```

## List of 2
## $ Male :List of 3
## ..$ Boomer : chr [1:10] "Theodore" "Bernard" "Gene" "Herbert" ...
## ..$ Genx : chr [1:10] "Marc" "Jamie" "Greg" "Darryl" ...
## ..$ Millennial: chr [1:10] "Zachary" "Dylan" "Christian" "Wesley" ...
## $ Female:List of 3
## ..$ Boomer : chr [1:10] "Gloria" "Joan" "Dorothy" "Shirley" ...
## ..$ Genx : chr [1:10] "Tracy" "Dawn" "Tina" "Tammy" ...
## ..$ Millennial: chr [1:10] "Samantha" "Alexis" "Brittany" "Lauren" ...

```

Dataframe creation

Imagine you are running an experiment with many manipulations. You have 14 levels in the variable “treatment” with levels a, b, c, etc. For each of those manipulations you have 3 submanipulations in a variable named “variation” with levels A, B, C. Then you have “gender” with levels M / F. Then you have “generation” with levels Boomer, GenX, Millenial. Then you will have 6 runs per each of these groups. In each set of 6 you will need to select a name without duplication from the appropriate set of names (from the last question). Create a data frame with columns treatment, variation, gender, generation, name and y that will store all the unique unit information in this experiment. Leave y empty because it will be measured as the experiment is executed.

```

n = 14 * 3 * 2 * 3 * 10
X = data.frame(
  treatment = rep(letters[1:14],n/14),
  variation = rep(LETTERS[1:3],n/3),
  Generations_Gender_list
)
#TO-DO

```

Packages

Install the package `pacman` using regular base R.

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

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.

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

- Create vector `v` consisting of all numbers from -100 to 100 and test using the second line of code su

```

v= seq(-100, 100)

#expect_equal(v, -100 : 101)

```

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

Test the `my_reverse` function from lab2 using the following code:

```
v = 1:100
#expect_equal(my_reverse(v), rev(v))
#expect_equal(my_reverse(c("A", "B", "C")), c("C", "B", "A"))
```

Multinomial Classification using KNN

Write a $(k=1)$ nearest neighbor algorithm using the Euclidean distance function. This is standard “Roxygen” format for documentation. Hopefully, we will get to packages at some point and we will go over this again. It is your job also to fill in this documentation.

```
#' Nearest neighbor classifier
#'
#' Classifying an observation based on its closest observation in the set.
#'
#' @param Xinput      A matrix of features for day
#' @param y_binary    A vector of training
#' @param Xtest       A test observation as a row vector.
#' @return            The predicted label for the test observation
nn_algorithm_predict = function(Xinput, y_binary, Xtest){
  n=nrow(Xinput)
  distances = array(NA,n)
  for (i in 1:n){
    distances [i] = sum((Xinput[i,]-Xtest)^2)
  }
  y_binary [which.min(distances)]
}
```

Write a few tests to ensure it actually works:

```
Xy = na.omit(MASS::biopsy) #The "breast cancer" data with all observations with missing values dropped
X = Xy[, 2 : 10] #V1, V2, ..., V9
y_binary = as.numeric(Xy$class == "malignant")
#pacman::p_load(class)
#y_1 = knn(X, c(4, 2, 1, 1, 2, 1, 2, 1, 1), y_binary, k = 1)
y_0 = nn_algorithm_predict (X,y_binary,c(4, 2, 1, 1, 2, 1, 2, 1, 1))
y_0
```

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

```
#y_1
```

We now add an argument `d` representing any legal distance function to the `nn_algorithm_predict` function. Update the implementation so it performs NN using that distance function. Set the default function to be the Euclidean distance in the original function. Also, alter the documentation in the appropriate places.

```
#' Nearest neighbor classifier
#'
#' Classifying an observation based on its closest observation in the set.
#'
#' @param Xinput      A matrix of features for day
#' @param y_binary    A vector of training
#' @param Xtest       A test observation as a row vector.
#' @param d           A distance function which takes 2 row vectors
#' @return            The predicted label for the test observation
nn_algorithm_predict = function(Xinput, y_binary, Xtest,d = function(v1,v2){sum((v1-v2)^2)}){

  n=nrow(Xinput)
  distances = array(NA,n)
  for (i in 1:n){
    distances [i] = d(Xinput[i,],Xtest)
  }
  y_binary [which.min(distances)]
}
```

For extra credit (unless you're a masters student), add an argument `k` to the `nn_algorithm_predict` function and update the implementation so it performs KNN. In the case of a tie, choose $\lfloor \hat{y} \rfloor$ randomly. Set the default `k` to be the square root of the size of $\lfloor \mathcal{D} \rfloor$ which is an empirical rule-of-thumb popularized by the “Pattern Classification” book by Duda, Hart and Stork (2007). Also,


```
mode = function(v){
  names(sort(table(v),decreasing = TRUE[1]))
}
```

- Fit a threshold model to `y` using the feature `Sepal.Length`. Write your own code to do this. What is the estimated value of the threshold parameter? Save the threshold value as `threshold`.

```
n = nrow(iris)
num_errors_by_parameter = matrix(NA, nrow = n, ncol = 2)
colnames(num_errors_by_parameter) = c("threshold_param", "num_errors")
y_logical = iris$Sepal.Length == "Yes"
for (i in 1 : n){
  threshold = iris$Sepal.Length[i]
  num_errors = sum((iris$Sepal.Length > threshold) != y_logical)
  num_errors_by_parameter[i, ] = c(threshold, num_errors)
}

threshold
```

```
## [1] 5.7
```

What is the total number of errors this model makes?

```
num_errors
```

```
## [1] 30
```

Does the threshold model's performance make sense given the following summaries:

```
threshold
```

```
## [1] 5.7
```

```
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
```

I would expect a lot more errors since it is a basic model.

Create the function `g` explicitly that can predict `y` from `x` being a new `Sepal.Length`.

```
g = function(x){
  if (x<= threshold) print ("Setosa")
  else print ("Versicolor")
}
```

Perceptron

You will code the “perceptron learning algorithm” for arbitrary number of features $\setminus(p)$. Take a look at the comments above the function. Respect the spec below:

```

#' Perceptron Algorithm
#'
#' A perceptron is an iterative algorithm that takes in a
#' linearly separable data set and returns a binary output based off an equation
#'
#' @param Xinput      Training data set matrix
#' @param y_binary    Binary training vector
#' @param MAX_ITER    Number of times the algorithm will run
#' @param w           A vector that will be created in the function
#'                   to store the training values (size p+1)
#'
#' @return            The computed final parameter (weight) as a vector of length p + 1

perceptron_learning_algorithm = function(Xinput, y_binary, MAX_ITER = 1000, w = NULL){

  p = ncol(Xinput)
  w = rep(0,p+1)
  Xinput = as.matrix(cbind(1,Xinput))

  for (t in 1 : MAX_ITER) {
    for (i in 1: nrow(Xinput)){
      x = Xinput[i,]
      y = y_binary[i]
      y_hat = (sum(x*w)> 0)
      for (j in 1:p){
        w[j] = w[j] + (y-y_hat) * x[j]
      }
    }
  }
  w
}

```

To understand what the algorithm is doing - linear “discrimination” between two response categories, we can draw a picture. First let’s make up some very simple training data \mathbb{D} .

```

Xy_simple = data.frame(
  response = factor(c(0, 0, 1, 1, 1)), #nominal
  first_feature = c(1, 1, 2, 3, 3, 4), #continuous
  second_feature = c(1, 2, 1, 3, 4, 3) #continuous
)

```

We haven’t spoken about visualization yet, but it is important we do some of it now. Thus, I will write this code for you and you will just run it. First we load the visualization library we’re going to use:

```
pacman::p_load(ggplot2)
```

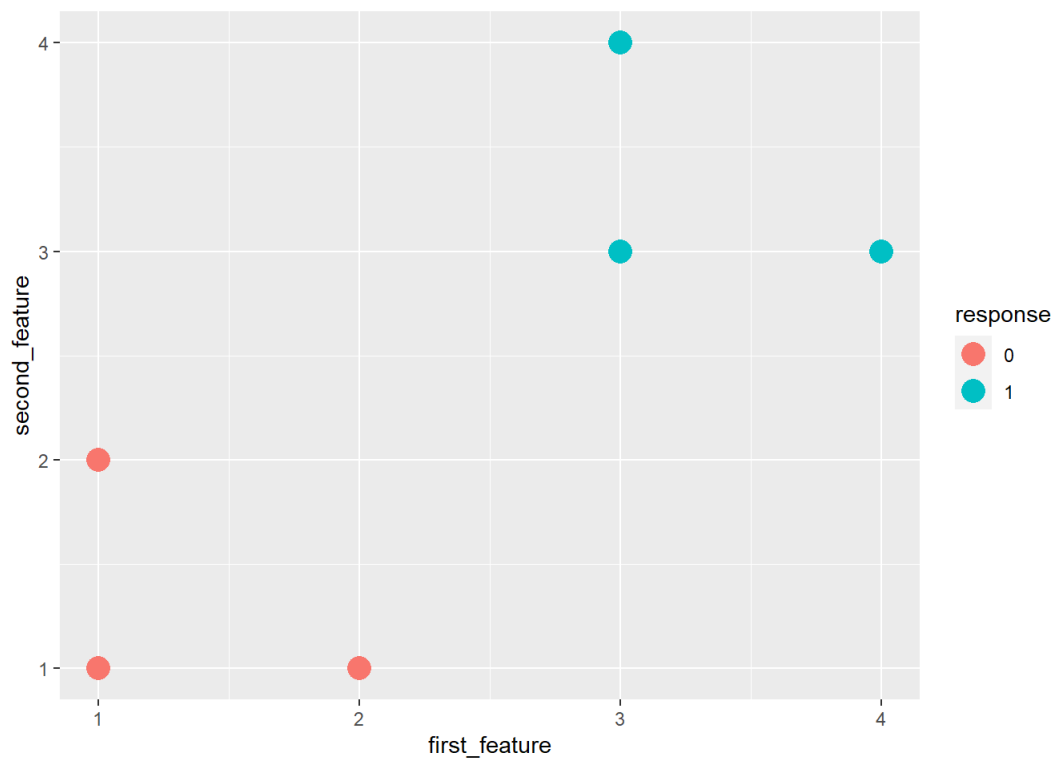
We are going to just get some plots and not talk about the code to generate them as we will have a whole unit on visualization using `ggplot2` in the future.

Let’s first plot y by the two features so the coordinate plane will be the two features and we use different colors to represent the third dimension, y .

```

simple_viz_obj = ggplot(Xy_simple, aes(x = first_feature, y = second_feature, color = response)) +
  geom_point(size = 5)
simple_viz_obj

```

TO-DO: We have two features and the data set is linearly separable

Now, let us run the algorithm and see what happens:

```
w_vec_simple_per = perceptron_learning_algorithm(
  cbind(Xy_simple$first_feature, Xy_simple$second_feature),
  as.numeric(Xy_simple$response == 1))
w_vec_simple_per
```

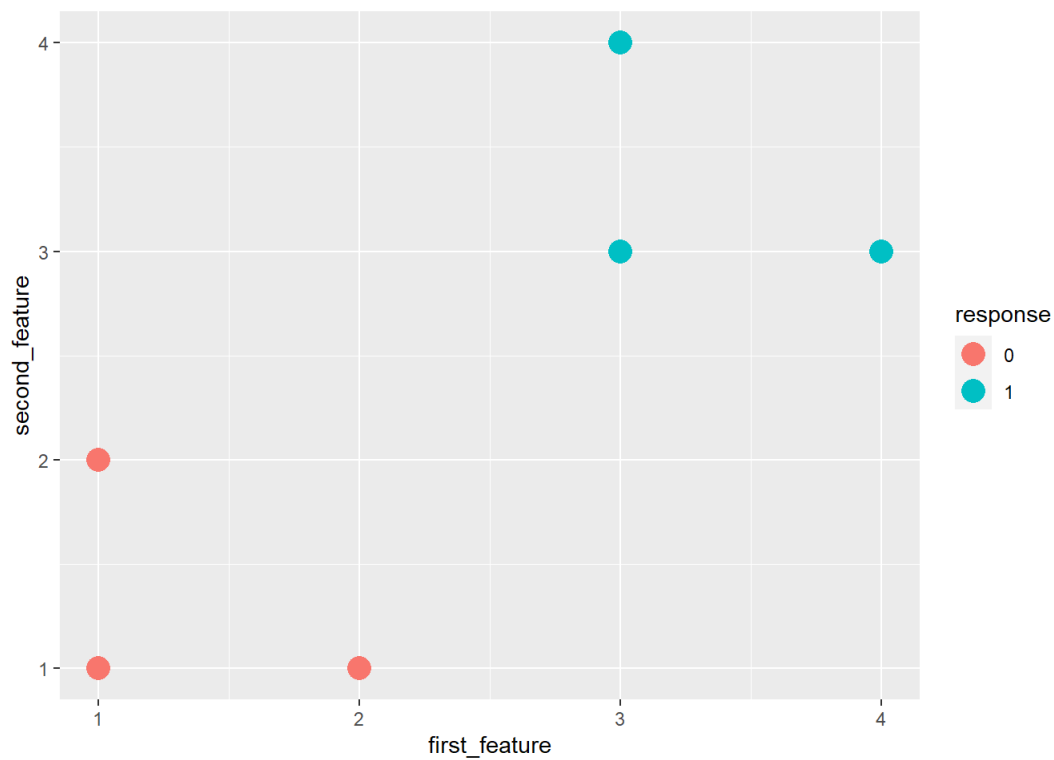
```
## [1] -2  1  0
```

Explain this output. What do the numbers mean? What is the intercept of this line and the slope? You will have to do some algebra.

$w_0 = -2$ $w_1 = 1$ $w_2 = 0$ The y-intercept is -2 with the slope being -1

```
simple_perceptron_line = geom_abline(
  intercept = -w_vec_simple_per[1] / w_vec_simple_per[3],
  slope = -w_vec_simple_per[2] / w_vec_simple_per[3],
  color = "orange")
simple_viz_obj + simple_perceptron_line
```

```
## Warning: Removed 1 rows containing missing values (geom_segment).
```



Explain this picture. Why is this line of separation not “satisfying” to you?

The line didn’t show up in the picture.

For extra credit, program the maximum-margin hyperplane perceptron that provides the best linear discrimination model for linearly separable data. Make sure you provide ROxygen documentation for this function.

```
#TO-DO
```