

```

---
  title: "Lab 2"
author: "Feiyi Ma"
output: pdf_document
date: "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.

```

```{r}
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)
```

```

* 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.

```

```{r}

flip_matrix = function(X,dim_to_rev=NULL){
 if(is.null(dim_to_rev)){
 dim_to_rev = ifelse(nrows(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)
X
flip_matrix(X,dim_to_rev = "cols")
```

```

* 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.

```
```{r}
arrays=list()
LETTERS
arrays[["A"]]=array(data=1:4, dim=c(2,2))
array(data=1:27,dim=c(3,3,3))
arrays
```
```

Run the following code:

```
```{r}
?lapply
lapply(arrays, object.size)
```
```

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?

```
```{r}
?object.size
object.size(arrays)
```
```

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

```
```{r}
#The object.size function is a length-one double value, an estimate of memory
allocation of the above value in one location. 232 bytes is smaller than 256
bytes that can be put into 512 bytes place.
```
```

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.

```
```{r}
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. "
?strsplit
x=sample(strsplit(lorem,split=",",),size=10,replace = TRUE, prob = NULL)
x
```
```

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.

```
```{r}
#HINT:
#strsplit("Theodore, Bernard, Gene, Herbert, Ray, Tom, Lee, Alfred, Leroy,
Eddie", split = ", ")[[1]]
?strsplit

m_b=strsplit("Theodore, Bernard, Gene, Herbert, Ray, Tom, Lee, Alfred, Leroy,
Eddie", split = ", ")[[1]]
m_g=strsplit("Marc, Jamie, Greg, Darryl, Tim, Dean, Jon, Chris, Troy, Jeff",
split = ", ")[[1]]
m_m=strsplit("Zachary, Dylan, Christian, Wesley, Seth, Austin, Gabriel, Evan,
Casey, Luis", split = ", ")[[1]]
f_b=strsplit("Gloria, Joan, Dorothy, Shirley, Betty, Dianne, Kay, Marjorie,
Lorraine, Mildred", split = ", ")[[1]]
f_g=strsplit("Tracy, Dawn, Tina, Tammy, Melinda, Tamara, Tracey, Colleen,
Sherri, Heidi", split = ", ")[[1]]
f_m=strsplit("Samantha, Alexis, Brittany, Lauren, Taylor, Bethany, Latoya,
Candice, Brittney, Cheyenne", split = ", ")[[1]]

b=list("M" = m_b, "F" = f_b)
g=list("M" = m_g, "F" = f_g)
m=list("M" = m_m, "F" = f_m)

data=list()
data[["Boomer"]]=b
data[["GenX"]]=g
data[["Millennial"]]=m

x=data
x
x$Boomer
x$GenX$F
```

## 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.

```
```{r}
n = 14 * 3 * 2 * 3 * 10
#X = data.frame(treatment = rep(NA,n),
?data.frame
treatment <- LETTERS[1:14]
variation<- letters[1:3]
t<-data.frame(treatment = rep(NA,n),1:14,sample(treatment,14,replace = TRUE))
v=data.frame(variation = rep(NA,n),1:3,sample(variation,3,replace = TRUE))
array(x,)
```

```
#TO-DO
```
```

Packages

Install the package `pacman` using regular base R.

```
```{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.

```
```{r}
pacman::p_load(testthat)
```
```

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

```
```{r}
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:

```
```{r}
v=1:100
expect_equal(my_reverse(v), rev(v))
expect_equal(my_reverse(c("A", "B", "C")), c("C", "B", "A"))
#if nothing, its equal
```
```

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.

```
```{r}
#' nearest Neighbor classifier
#'
#' classify an observation based on label of the closest observation in the set
of traning observations
#'
#' @param Xinput A matrix of features for training data observation
#' @param y_binary The vector of training data label
#' @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) # ith row vector distance
difference
 }
 y_binary[which.min(distances)]
}
```
```

Write a few tests to ensure it actually works:

```
```{r}
?sample
x=sample(c(1:10),100,replace = TRUE,prob=NULL)
a=matrix(data=x,nrow=10,ncol=10,byrow=FALSE,dimnames = NULL)
y=matrix(data=NA,nrow=10,ncol = 10,byrow = TRUE)
t=sample(c(1:10),1)
nn_algorithm_predict(a,y, t)
?which.min
```
```

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.

```

```{r}
#' nearest Neighbor classifier
#'
#' classify an observation based on label of the closest observation in the set
of training observations
#'
#' @param Xinput A matrix of features for training data observation
#' @param y_binary The vector of training data label
#' @param xtest A test observation as a row vector
#' @param d A distance function which take inputs to two different
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) # ith row vector distance difference
 }
 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 \hat{y} randomly. Set the default `k` to be the square root of the size of \mathcal{D} which is an empirical rule-of-thumb popularized by the "Pattern Classification" book by Duda, Hart and Stork (2007). Also, alter the documentation in the appropriate places.

```

```{r}
#TO-DO for the 650 students but extra credit for undergrads
#' nearest Neighbor classifier
#'
#' classify an observation based on label of the closest observation in the set
of training observations
#'
#' @param Xinput A matrix of features for training data observation
#' @param y_binary The vector of training data label
#' @param xtest A test observation as a row vector
#' @param d A distance function which take inputs to two different
row vectors.
#' @param k closed people of the k sample, find the result of prob
#' @return The predicted label for the test observation
nn_algorithm_predict = function(Xinput, y_binary, xtest, d=function(v1,v2)
{sum((v1-v2)^2)}, k) {
 n=nrow(Xinput)
 distances=array(NA,n)

```

```

 for(i in 1:n){
 distances[i]=d(xinput[i,],xtext) # ith row vector distance difference
 }
 y_binary[which.min(distances)]
 }
}

```

## ## Basic Binary Classification Modeling

\* Load the famous `iris` data frame into the namespace. Provide a summary of the columns using the `skim` function in package `skimr` and write a few descriptive sentences about the distributions using the code below and in English.

```

```{r}
data(iris)
pacman::p_load(skimr)
skim(iris)
#summary(iris)
```

```

TO-DO: classify iris flowers based on petal length, width and sepal length and width. The above chart show how the iris flower classification based on these data. The average sepal length is 5.47, the average sepal width is 3.10, the average petal length is 2.86, the average petal width is 0.786.

The outcome / label / response is `Species`. This is what we will be trying to predict. However, we only care about binary classification between "setosa" and "versicolor" for the purposes of this exercise. Thus the first order of business is to drop one class. Let's drop the data for the level "virginica" from the data frame.

```

```{r}
#when sepal not equal to petal
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".

```

```{r}
y=as.integer(iris$Species == "setosa")
```

```

\* Write a function `mode` returning the sample mode.

```

```{r}
mode=function(v) {
  names(sort(table(v), decreasing=TRUE)[1])
}
mode(iris3)
```

```

\* 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`.

```
```{r}
x=as.matrix(iris$Sepal.Length)
n = nrow(x)
num_errors_by_parameter = matrix(NA, nrow = n, ncol = 2)
colnames(num_errors_by_parameter) = c("threshold_param", "num_errors")

for (i in 1 : n){
  threshold = x[i,]
  num_errors = sum((x > threshold) != y)
  num_errors_by_parameter[i, ] = c(threshold, num_errors)
}

num_errors_by_parameter[order(num_errors_by_parameter[, "num_errors"]), ]
#now grab the smallest num errors
best_row = order(num_errors_by_parameter[, "num_errors"])[1]
shreshold = c(num_errors_by_parameter[best_row, "threshold_param"], use.names
= FALSE)
shreshold
```
```

What is the total number of errors this model makes?

```
```{r}
n_errors=sum(num_errors_by_parameter[,2])
n_errors
```
```

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

```
```{r}
threshold
summary(iris[iris$Species == "setosa", "Sepal.Length"])
summary(iris[iris$Species == "versicolor", "Sepal.Length"])
```
```

TO-DO: threshold model result is 5.9, which is the very closed to the mean using summary function, it is in the reasonable range of error.

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

```
```{r}
g = function(x){
  ifelse(y>shreshold,1,0)
}

```
```



## Perceptron

You will code the "perceptron learning algorithm" for arbitrary number of features  $p$ . Take a look at the comments above the function. Respect the spec below:

```
```{r}
#' TO-DO: perceptron learning algorithm
#'
#' TO-DO: Explain what this function does in a few sentences
#'
#' @param Xinput      A matrix of features for training data observation
#' @param y_binary    Real value from the data
#' @param MAX_ITER    number of sample observations
#' @param w           random vector
#'
#' @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) {
  #TO-DO

  w_vec = rep(0, dim(x)) #intialize a n-dim vector
  X1 = as.matrix(cbind(1, x[, 1, drop = FALSE]))
  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_binary[i]

      w_vec[1] = w_vec[1] + (y_i - yhat_i) * x_i[1]
      w_vec[2] = w_vec[2] + (y_i - yhat_i) * x_i[2]

    }
  }
  #dont know how to change the loop to vector loop
  w_vec
}
```
```

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  $\mathbb{D}$ .

```
```{r}
Xy_simple = data.frame(
  response = factor(c(0, 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:

```
```{r}
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,  $Y$ .

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

TO-DO: it shows the sample data set in 3D. response is the color, first\_feature is the x\_axis and second\_feature is the y\_axis. the vector value is stored in the three data columns. EG. red (1,1) is saved as response =0, first\_feature =1, second\_feature = 1.

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

```
```{r}
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
w_norm = sqrt(sum(w_vec_simple_per^2))

#perceptron_learning_algorithm code not sure how to generalized to n
dimension.
```
```

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.

first number is the b term, second is  $w_1$ , last is  $w_2$

```
```{r}
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
```
```

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

TO-DO

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.

```
```{r}  
#TO-DO  
```
```