

# 36350-A HW7

Joong Ho Choi

TOTAL POINTS

**110 / 120**

## QUESTION 1

1 Q1 16 / 20

- 0 pts Correct
- 4 pts does not display histogram with the estimated density
- 2 pts x values not between 0 and 4
- ✓ - 4 pts incorrect estimation (too far from 0.666)
- 4 pts question links incorrectly
- 20 pts blank

## QUESTION 2

2 Q2 20 / 20

- ✓ - 0 pts Correct
- 2 pts incorrect lower bounds
- 2 pts incorrect upper bounds
- 4 pts function entered incorrectly
- 3 pts incorrect answer
- 4 pts question links incorrectly
- 20 pts blank

## QUESTION 3

3 Q3 20 / 20

- ✓ - 0 pts Correct
- 3 pts does not sample 100,000 integers between 1 and 100 with replacement
- 2 pts incorrect matrix dimension (should have 1000 rows and 100 columns)
- 4 pts incorrect # of unique values in each row
- 1 pts does not save # of unique values into a vector
- 3 pts no histogram
- 1 pts no appropriate labels
- 2 pts missing / incorrect mean
- 2 pts missing / incorrect s.d.

- 4 pts question links incorrectly
- 20 pts blank

## QUESTION 4

4 Q4 20 / 20

- ✓ - 0 pts Correct
- 3 pts does not display scatter plot with fitted curve
- 4 pts does not use MSE / Wrong MSE
- 2 pts does not use "BFGS" in the optim function
- 4 pts curve not fit data well / wrong gradient
- 4 pts incorrect initial estimated parameter values
- 4 pts question links incorrectly
- 20 pts blank

## QUESTION 5

5 Q5 18 / 20

- 0 pts Correct
- ✓ - 2 pts does not return null when input is not character type
- 2 pts does not return null when input is not a single string
- 2 pts does not convert input string to lowercase
- 2 pts does not use permn function
- 2 pts does not show output for "rneup"
- 2 pts does not show output for "srsets"
- 2 pts does not show output for "lypelu"
- 4 pts does not return unique words in GradyAugmented
- 4 pts question links incorrectly
- 2 pts question partially linked
- 20 pts blank

## QUESTION 6

6 Q6 16 / 20

- 0 pts Correct
- 2 pts does not generate 10,000 permutations

- **6 pts** incorrect % of differences greater than the observed difference
- **6 pts** incorrect conclusion about null hypothesis
- **4 pts** sampled with replacement
- **6 pts** Did not sample from  $x$
- **6 pts** Did not sample the correct number of values

(20)

- **2 pts** did not divide by  $k = 10000$
- **4 pts** Does not take proportion to generate p-value

- **6 pts** modifies  $y$
- **20 pts** blank
- **4 pts** question links incorrectly

**- 4 Point adjustment**

💬 please print p-value

QUESTION 7

7 Late Penalty 0 / 0

✓ - **0 pts** Correct

# HW: Week 7

36-350 – Statistical Computing

Week 7 – Spring 2021

Name: Joong Ho Choi

Andrew ID: joonghoc

You must submit **your own** HW as a PDF file on Gradescope.

---

## HW Length Cap Instructions

- If the question requires you to print a data frame in your solution e.g. `q1_out_df`, you must first apply `head(q1_out_df, 30)` and `dim(q1_out_df)` in the final knitted pdf output for such a data frame.
- Please note that this only applies if you are knitting the `Rmd` to a `pdf`, for Gradescope submission purposes.
- If you are using the data frame output for visualization purposes (for example), use the entire data frame in your exploration
- The **maximum allowable length** of knitted pdf HW submission is **30 pages**. Submissions exceeding this length *will not be graded* by the TAs. All pages must be tagged as usual for the required questions per the usual policy
- For any concerns about HW length for submission, please reach out on Piazza during office hours

## Question 1

(20 points)

Let's say you have data sampled randomly from an unknown distribution. How can you estimate the probability that the next datum  $x_o$  that you observe will lie between the values  $x = a$  and  $x = b$ ?

One technique is to guess the distribution from a set of named ones, go through the optimization process, and then perform the necessary integral, by hand or with `integrate()`. But if you cannot guess the distribution, you could also use density estimation, i.e., nonparametrically estimate the probability density function using a kernel density estimator. Given the estimated pdf, you would then compute the required integral.

Below is a histogram for a data sample of size 100. Use the `density()` function to estimate the underlying pdf. By default, the weighting kernel is "`gaussian`" and the width of the kernel is computed via a "plug-in rule"; you do not need to change these default values. Overlay the estimated density function onto the histogram.

Then compute the probability between  $x = 0$  and  $x = 4$ , using the output from `density()`. You can code your own integrator, based on summations (trapezoid rule!), or utilize an appropriate function in the `sfsmisc` package that will integrate a function that is defined as a series of  $(x, y)$  pairs and not as a parametrized function. Note that my answer is approximately 0.666; you do not need to match this *exactly*, but your own answer should be pretty close. (Remember: in the end, this is only an estimate of the probability. Changing

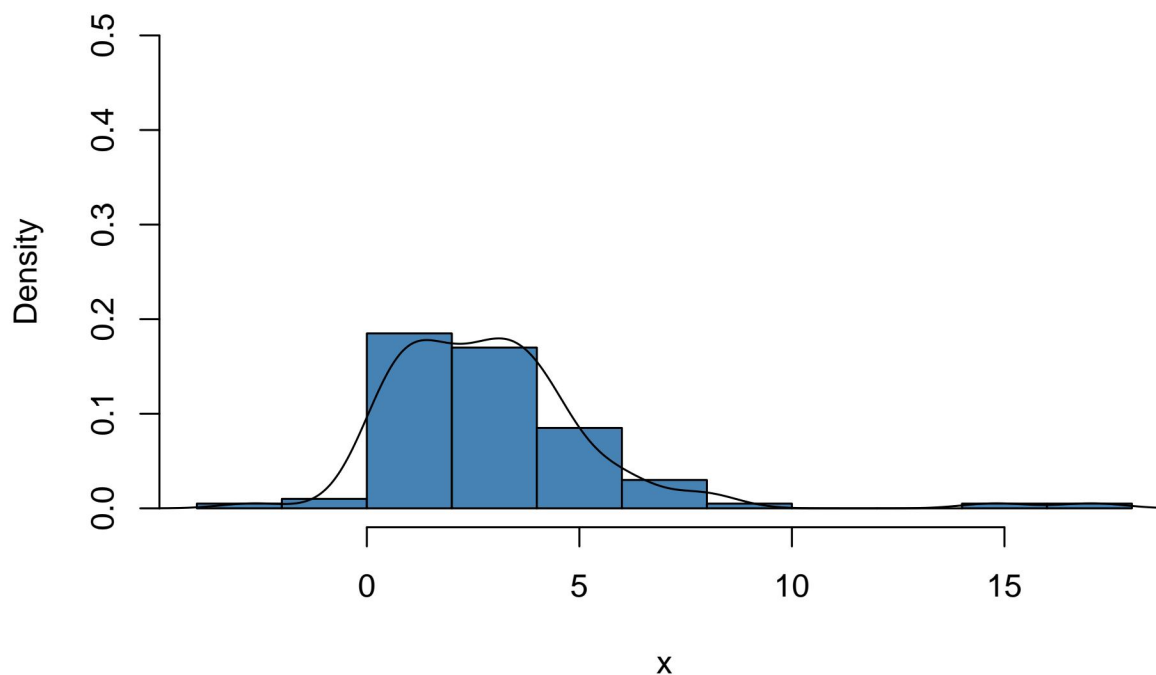
the width and/or form of the kernel will change your answer. Also, a kernel density estimate smooths data, so you can never actually recover the true underlying pdf anyway, just an approximation of it that is going to be slightly too wide. Still, this whole process beats having to fall back on Tchebysheff's theorem.)

```
if ( require(sfsmisc) == FALSE ) {  
  install.packages("sfsmisc",repos="https://cloud.r-project.org")  
  library(sfsmisc)  
}
```

```
## Loading required package: sfsmisc
```

```
set.seed(303)  
x = rgamma(100,2.5,scale=1.25)  
s = sample(c(-1,1),100,replace=TRUE)  
x = x + s*rexp(100,rate=3/abs(x))  
#hist(x,prob=TRUE,ylim=c(0,0.5),col="steelblue")  
hist(x,prob=TRUE,ylim=c(0,0.5),col="steelblue")  
lines(density(x))
```

## Histogram of x



```
d_fun <- ecdf (x)  
d_fun(4) -d_fun(0) #Area under curve <= 4 - area under curve <=0
```

```
## [1] 0.71
```

1 Q1 16 / 20

- 0 pts Correct
- 4 pts does not display histogram with the estimated density
- 2 pts x values not between 0 and 4
- ✓ - 4 pts incorrect estimation (too far from 0.666)
- 4 pts question links incorrectly
- 20 pts blank

## Question 2

(20 points)

The `integrate()` function is used to perform univariate integrals of parametrized functions. What if you have a multi-dimensional integral instead, like

$$\int_0^{\pi/2} \int_0^{\pi/4} \sin(2x + y) \cos(x + 2y) dx dy ?$$

One option is to use the `cubature` package (“adaptive multivariate integration over hypercubes”). Install and use a function or functions of the `cubature` package to compute the integral given above. (Note: this package is useful for *hypercubes*; if the bounds of integration are variable [like when you dealt with triangular regions of integration in 225], then the functions of the `cubature` package are not appropriate. At that point, you would start moving towards, e.g., Monte Carlo integration.)

```
if ( require(cubature) == FALSE ) {  
  install.packages("cubature",repos="https://cloud.r-project.org")  
  library(cubature)  
}
```

```
## Loading required package: cubature
```

```
f<-function(x){  
  sin(2*x[1]+x[2])*cos(x[1]+2*x[2])  
}  
  
hcubature(f,lower=c(0,0),upper=c(pi/4,pi/2))
```

```
## $integral  
## [1] -0.2626623  
##  
## $error  
## [1] 2.421615e-06  
##  
## $functionEvaluations  
## [1] 697  
##  
## $returnCode  
## [1] 0
```

## Question 3

(20 points)

If you sample 100 integers between 1 and 100 inclusive with replacement, how many unique integers are selected, on average? To answer this, do the following:

- (1) sample 100,000 integers between 1 and 100 inclusive with replacement;
- (2) place these samples into a  $1000 \times 100$  matrix (so you’ve effectively simulated 1000 datasets with 100 samples each;
- (3) determine the number of unique values in each row and save the result to a vector;

## 2 Q2 20 / 20

✓ - 0 pts Correct

- 2 pts incorrect lower bounds
- 2 pts incorrect upper bounds
- 4 pts function entered incorrectly
- 3 pts incorrect answer
- 4 pts question links incorrectly
- 20 pts blank

## Question 2

(20 points)

The `integrate()` function is used to perform univariate integrals of parametrized functions. What if you have a multi-dimensional integral instead, like

$$\int_0^{\pi/2} \int_0^{\pi/4} \sin(2x + y) \cos(x + 2y) dx dy ?$$

One option is to use the `cubature` package (“adaptive multivariate integration over hypercubes”). Install and use a function or functions of the `cubature` package to compute the integral given above. (Note: this package is useful for *hypercubes*; if the bounds of integration are variable [like when you dealt with triangular regions of integration in 225], then the functions of the `cubature` package are not appropriate. At that point, you would start moving towards, e.g., Monte Carlo integration.)

```
if ( require(cubature) == FALSE ) {  
  install.packages("cubature",repos="https://cloud.r-project.org")  
  library(cubature)  
}
```

```
## Loading required package: cubature
```

```
f<-function(x){  
  sin(2*x[1]+x[2])*cos(x[1]+2*x[2])  
}  
  
hcubature(f,lower=c(0,0),upper=c(pi/4,pi/2))
```

```
## $integral  
## [1] -0.2626623  
##  
## $error  
## [1] 2.421615e-06  
##  
## $functionEvaluations  
## [1] 697  
##  
## $returnCode  
## [1] 0
```

## Question 3

(20 points)

If you sample 100 integers between 1 and 100 inclusive with replacement, how many unique integers are selected, on average? To answer this, do the following:

- (1) sample 100,000 integers between 1 and 100 inclusive with replacement;
- (2) place these samples into a  $1000 \times 100$  matrix (so you’ve effectively simulated 1000 datasets with 100 samples each;
- (3) determine the number of unique values in each row and save the result to a vector;



- (4) histogram the result, properly labeled (base R or `ggplot()`); and
- (5) compute the sample mean and sample standard deviation of vector generated in step (3). If the mean is much different from 63.397, you've made a misstep. (Now, why are you doing this? Because the result tells you that in any one bootstrap sample of your data, there is a roughly 63% chance that a particular datum will be sampled one or more times. Or, alternatively, that roughly 37% of your data will not be sampled. Google, e.g., "out-of-bag error".)

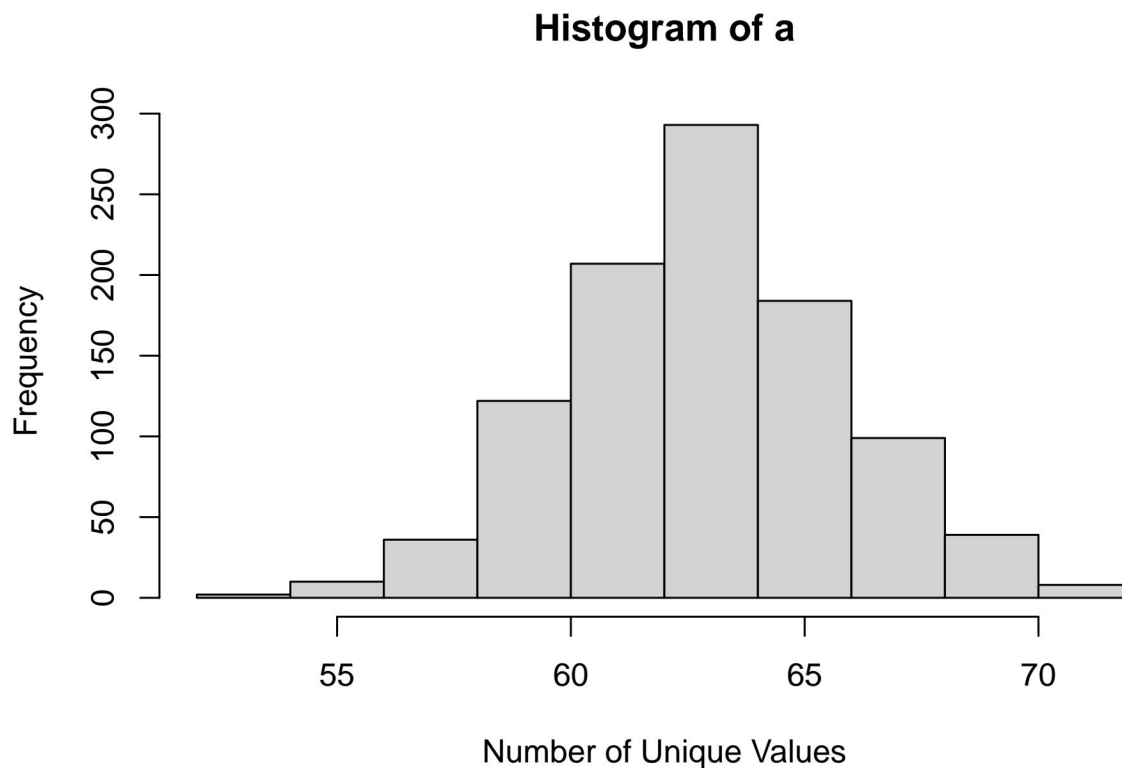
```
#library(ggplot2)
k<-sample(c(1:100), 100000, replace = TRUE)

matrix_x<-matrix(unlist(k),nrow=1000,ncol=100)

a<-rep(NA,1000)  #vectors of step(3)

for (ii in 1:1000){
  a[ii]=length((unique(matrix_x[ii,])))
}

hist(a,xlab="Number of Unique Values")
```



```
mean(a)
```

```
## [1] 63.355
```

```
sd(a)
```

```
## [1] 2.988965
```

## Question 4

(20 points)

Below we load in a set of  $(x, y)$  data pairs:

```
load(url("http://www.stat.cmu.edu/~pfreeman/HW_07_Q7.Rdata"))
```

Fit a quadratic function  $f(x) = ax^2 + bx + c$  to these data using `optim()`, with a gradient function specified. Use `method = BFGS`, and instead of the negative log-likelihood, minimize the residual sum of squares (i.e., minimize the squared distance from the model to the data).

Plot the data and overplot the best-fit line. Note that if that line does not match the data well, you might have a situation in which your initial guess of the parameter values was incorrect... so try fitting again with a new guess. (Hint: remember that the gradient is the gradient of the fit metric, not the derivative of  $f(x)$ , and that the gradient function returns a vector of length three: the partial derivative of the fit metric with respect to  $a$ , then with respect to  $b$ , then with respect to  $c$ .)

You'll need to derive the gradient vector by hand. Ask us sooner rather than later how to do this if you've forgotten... but you should have covered this in 21-259 or the equivalent.)

```
dat=data.frame(x=x,y=y)
f<-function(x,my.par){
  my.par[[1]]*x^2+my.par[[2]]*x+my.par[[3]]
}

SSE<-function(my.par){sum((y-f(x,my.par))^2)}

fx=expression((y-(a*x^2+b*x+c)))
(df1=D(fx,'a')) # -x^2

## -x^2

(df2=D(fx,'b')) # -x

## -x

(df3=D(fx,'c')) # -1

## -1

SSE.fit<-function(my.par){
  c(sum(-x^2),sum(-x),sum(-1))
}

#(result <- optim(par = c(2,2,2), fn = SSE,gr=SSE.fit,method="BFGS"))

#(result <- optim(par = c(0.2531734,-1.3146661,-0.6579489), fn = SSE,gr=SSE.fit,method="BFGS"))
(result <- optim(par = c(2,2,2), fn = SSE,method="BFGS"))
```

### 3 Q3 20 / 20

✓ - 0 pts Correct

- 3 pts does not sample 100,000 integers between 1 and 100 with replacement
- 2 pts incorrect matrix dimension (should have 1000 rows and 100 columns)
- 4 pts incorrect # of unique values in each row
- 1 pts does not save # of unique values into a vector
- 3 pts no histogram
- 1 pts no appropriate labels
- 2 pts missing / incorrect mean
- 2 pts missing / incorrect s.d.
- 4 pts question links incorrectly
- 20 pts blank

```
sd(a)
```

```
## [1] 2.988965
```

## Question 4

(20 points)

Below we load in a set of  $(x, y)$  data pairs:

```
load(url("http://www.stat.cmu.edu/~pfreeman/HW_07_Q7.Rdata"))
```

Fit a quadratic function  $f(x) = ax^2 + bx + c$  to these data using `optim()`, with a gradient function specified. Use `method = BFGS`, and instead of the negative log-likelihood, minimize the residual sum of squares (i.e., minimize the squared distance from the model to the data).

Plot the data and overplot the best-fit line. Note that if that line does not match the data well, you might have a situation in which your initial guess of the parameter values was incorrect... so try fitting again with a new guess. (Hint: remember that the gradient is the gradient of the fit metric, not the derivative of  $f(x)$ , and that the gradient function returns a vector of length three: the partial derivative of the fit metric with respect to  $a$ , then with respect to  $b$ , then with respect to  $c$ .)

You'll need to derive the gradient vector by hand. Ask us sooner rather than later how to do this if you've forgotten... but you should have covered this in 21-259 or the equivalent.)

```
dat=data.frame(x=x,y=y)
f<-function(x,my.par){
  my.par[[1]]*x^2+my.par[[2]]*x+my.par[[3]]
}

SSE<-function(my.par){sum((y-f(x,my.par))^2)}

fx=expression((y-(a*x^2+b*x+c)))
(df1=D(fx,'a')) # -x^2

## -x^2

(df2=D(fx,'b')) # -x

## -x

(df3=D(fx,'c')) # -1

## -1

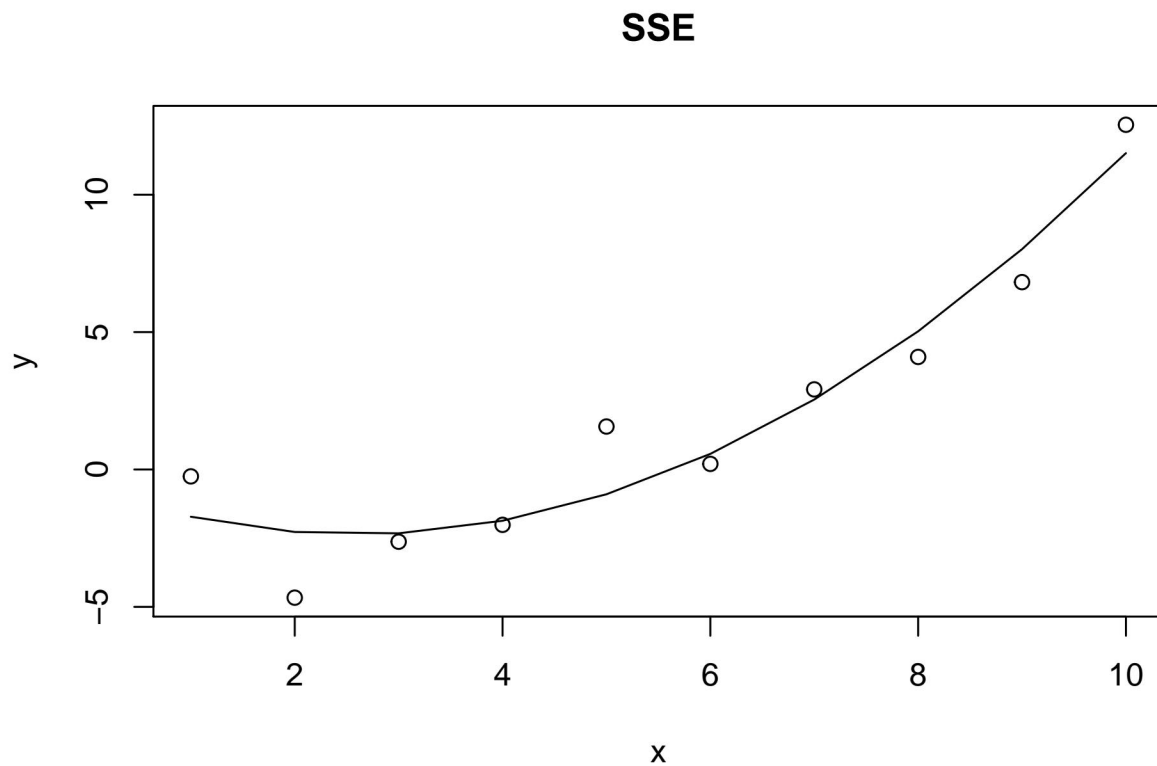
SSE.fit<-function(my.par){
  c(sum(-x^2),sum(-x),sum(-1))
}

#(result <- optim(par = c(2,2,2), fn = SSE,gr=SSE.fit,method="BFGS"))

#(result <- optim(par = c(0.2531734,-1.3146661,-0.6579489), fn = SSE,gr=SSE.fit,method="BFGS"))
(result <- optim(par = c(2,2,2), fn = SSE,method="BFGS"))
```

```
## $par
## [1] 0.2531734 -1.3146661 -0.6579489
##
## $value
## [1] 17.70249
##
## $counts
## function gradient
##      27      7
##
## $convergence
## [1] 0
##
## $message
## NULL
```

```
plot(y ~ x, data=dat, main='SSE')
lines(result$par[[1]]*x^2+result$par[[2]]*x+result$par[[3]])
```



## Question 5

*(20 points)*

Jumble, the scrambled-word game that is syndicated in many newspapers, can be hard. If we can have the computer play Jumble for us, well... it would be easier. Maybe not as fulfilling, but certainly easier.

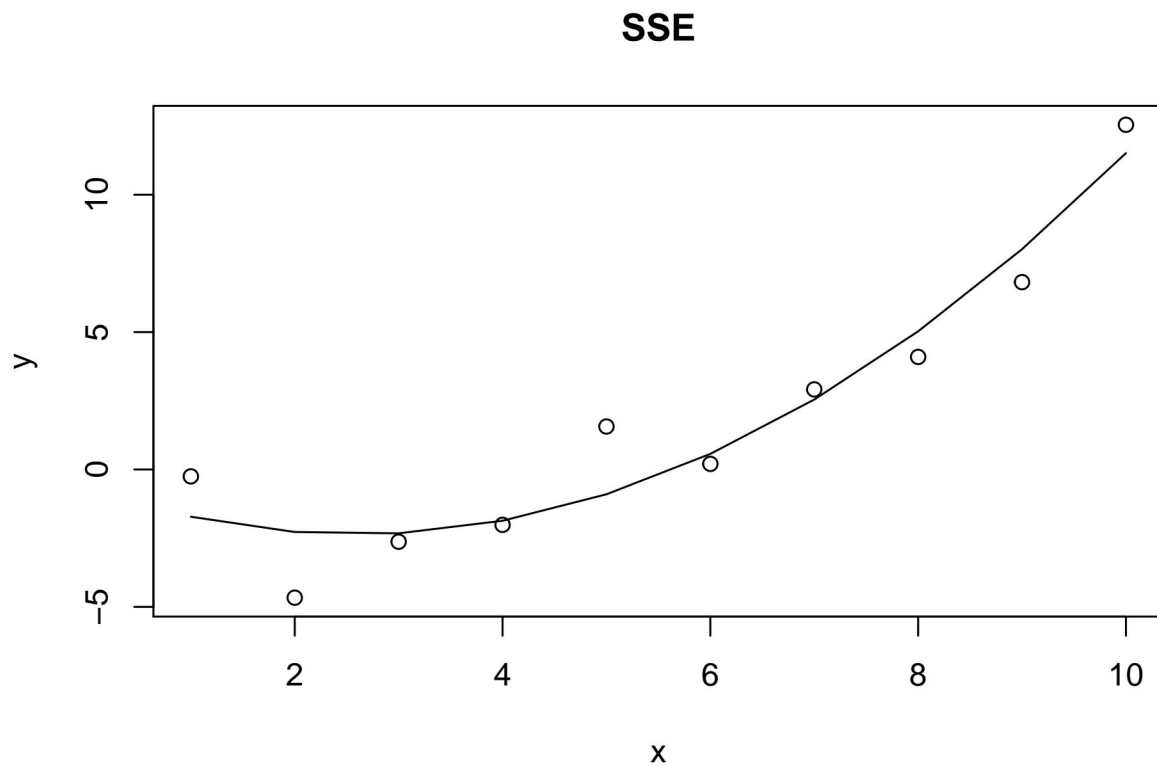
#### 4 Q4 20 / 20

✓ - 0 pts Correct

- 3 pts does not display scatter plot with fitted curve
- 4 pts does not use MSE / Wrong MSE
- 2 pts does not use "BFGS" in the optim function
- 4 pts curve not fit data well / wrong gradient
- 4 pts incorrect initial estimated parameter values
- 4 pts question links incorrectly
- 20 pts blank

```
## $par
## [1] 0.2531734 -1.3146661 -0.6579489
##
## $value
## [1] 17.70249
##
## $counts
## function gradient
##      27      7
##
## $convergence
## [1] 0
##
## $message
## NULL
```

```
plot(y ~ x, data=dat, main='SSE')
lines(result$par[[1]]*x^2+result$par[[2]]*x+result$par[[3]])
```



## Question 5

*(20 points)*

Jumble, the scrambled-word game that is syndicated in many newspapers, can be hard. If we can have the computer play Jumble for us, well... it would be easier. Maybe not as fulfilling, but certainly easier.

Construct a function called `matchDict` that utilizes the `permn()` function of the `combinat` package and the 120,000+ length character vector `GradyAugmented` in the `qdapDictionaries` package so as to return possible unscrambled solutions given a scrambled input.

`permn()` will output a list of all permutations of the letters in the input word; you need to combine these letters back into candidate words and determine which of the candidate words is in the `GradyAugmented` vector.

(`permn()` is a useful tool for generating permutations, hence the idea for this exercise.) For instance, if you input “nidkr”, you should get “drink” as your only output. Show output for “rneup”, “srsets”, and “LYPELU”.

Note: if your input is not a single string, have your function return `NULL`, and if your input is an upper-case string, convert it to lower-case: the `GradyAugmented` vector consists of lower-case strings.

Also, return only unique instances of words. Hint: if you need to paste letters together, use `paste()` and look closely at its arguments to make sure you set the right argument to the right value.

```
if ( require(combinat) == FALSE ) {
  install.packages("combinat",repos="https://cloud.r-project.org")
  library(combinat)
}
```

```
## Loading required package: combinat
```

```
##
```

```
## Attaching package: 'combinat'
```

```
## The following object is masked from 'package:utils':
```

```
##
```

```
##      combn
```

```
if ( require(qdapDictionaries) == FALSE ) {
  install.packages("qdapDictionaries",repos="https://cloud.r-project.org")
  library(qdapDictionaries)
}
```

```
## Loading required package: qdapDictionaries
```

```
matchDict<-function(word){
  if (length(word)!=1){
    return (NULL)
  }
  word<-tolower(word)
  w<-as.character(word)
  w<-unlist(strsplit(w,split = ""))
  choices=permn(w) #choices=list
  b<-lapply(choices,function(x) paste(x,sep="",collapse=""))
  i<-which(b %in% GradyAugmented)
  return (unique(b[i]))
}
matchDict("nidkr")
```

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

```
## [1] "drink"
```



```
matchDict("rneup")
```

```
## [[1]]  
## [1] "prune"
```

```
matchDict("srsets")
```

```
## [[1]]  
## [1] "stress"
```

```
matchDict("LYPELU")
```

```
## [[1]]  
## [1] "pulley"
```

## Question 6

*(20 points)*

The bootstrap is one type of what are called “resampling tests.” Another type of resampling test is the permutation test. Assume you observe a series of  $(x, y)$  pairs; for instance,  $x$  might be a factor variable with two levels representing treatment groups, and  $y$  might be the observed responses. In this context, you might be interested in seeing if the difference in the mean values of  $y$  for each treatment group is either significantly different from zero (a two-sided hypothesis test) or significantly less than or greater than zero (a one-sided test).

In a permutation test, you estimate the p-value by randomly shuffling the  $x$  vector (while leaving the  $y$  vector intact!), then computing the difference, and repeating until you build up a vector of differences, and seeing, e.g., how many of the differences generated via permutation are greater than the actual difference you observe.

In the code chunk below,  $x$  and  $y$  are the observed data. Generate the observed difference between the means for the groups  $x = 0$  and  $x = 1$ .

Then code a permutation test. Note that you cannot use `permn()` here to generate all the possible permutations for the  $x$  vector, as that would be  $\sim 10^{18}$  permutations. So: generate 10,000 permutations, save the differences for each, and see how many of the 10,000 differences are greater than what you observe. In the end, do you reject or fail to reject the null hypothesis that the actual difference in population means is zero?

```
set.seed(1002)  
x = c(rep(0,10),rep(1,10))  
y = x+rnorm(20,sd=1.5)  
observed.diff = diff(tapply(y,x,mean))  
  
res<-replicate(10000,diff(tapply(y,sample(x,length(x),FALSE),mean)))  
a<-length(which(res>observed.diff))  
p=a/10000  
if (p<0.05) cat("reject the null hypothesis")
```

```
## reject the null hypothesis
```

5 Q5 18 / 20

- 0 pts Correct

✓ - 2 pts does not return null when input is not character type

- 2 pts does not return null when input is not a single string

- 2 pts does not convert input string to lowercase

- 2 pts does not use permn function

- 2 pts does not show output for "rneup"

- 2 pts does not show output for "srsets"

- 2 pts does not show output for "lypelu"

- 4 pts does not return unique words in GradyAugmented

- 4 pts question links incorrectly

- 2 pts question partially linked

- 20 pts blank

```
matchDict("rneup")
```

```
## [[1]]  
## [1] "prune"
```

```
matchDict("srsets")
```

```
## [[1]]  
## [1] "stress"
```

```
matchDict("LYPELU")
```

```
## [[1]]  
## [1] "pulley"
```

## Question 6

*(20 points)*

The bootstrap is one type of what are called “resampling tests.” Another type of resampling test is the permutation test. Assume you observe a series of  $(x, y)$  pairs; for instance,  $x$  might be a factor variable with two levels representing treatment groups, and  $y$  might be the observed responses. In this context, you might be interested in seeing if the difference in the mean values of  $y$  for each treatment group is either significantly different from zero (a two-sided hypothesis test) or significantly less than or greater than zero (a one-sided test).

In a permutation test, you estimate the p-value by randomly shuffling the  $x$  vector (while leaving the  $y$  vector intact!), then computing the difference, and repeating until you build up a vector of differences, and seeing, e.g., how many of the differences generated via permutation are greater than the actual difference you observe.

In the code chunk below,  $x$  and  $y$  are the observed data. Generate the observed difference between the means for the groups  $x = 0$  and  $x = 1$ .

Then code a permutation test. Note that you cannot use `permn()` here to generate all the possible permutations for the  $x$  vector, as that would be  $\sim 10^{18}$  permutations. So: generate 10,000 permutations, save the differences for each, and see how many of the 10,000 differences are greater than what you observe. In the end, do you reject or fail to reject the null hypothesis that the actual difference in population means is zero?

```
set.seed(1002)  
x = c(rep(0,10),rep(1,10))  
y = x+rnorm(20,sd=1.5)  
observed.diff = diff(tapply(y,x,mean))  
  
res<-replicate(10000,diff(tapply(y,sample(x,length(x),FALSE),mean)))  
a<-length(which(res>observed.diff))  
p=a/10000  
if (p<0.05) cat("reject the null hypothesis")
```

```
## reject the null hypothesis
```

```
if (p>0.05) cat("Fail to reject the null hypothesis")
```

FILL ME IN: be sure to break lines so that all your text appears in the PDF file.

6 Q6 16 / 20

- 0 pts Correct
- 2 pts does not generate 10,000 permutations
- 6 pts incorrect % of differences greater than the observed difference
- 6 pts incorrect conclusion about null hypothesis
- 4 pts sampled with replacement
- 6 pts Did not sample from x
- 6 pts Did not sample the correct number of values (20)
- 2 pts did not divide by  $k = 10000$
- 4 pts Does not take proportion to generate p-value
- 6 pts modifies y
- 20 pts blank
- 4 pts question links incorrectly
- 4 Point adjustment
  - 🗨 please print p-value

7 Late Penalty 0 / 0

✓ - 0 pts Correct