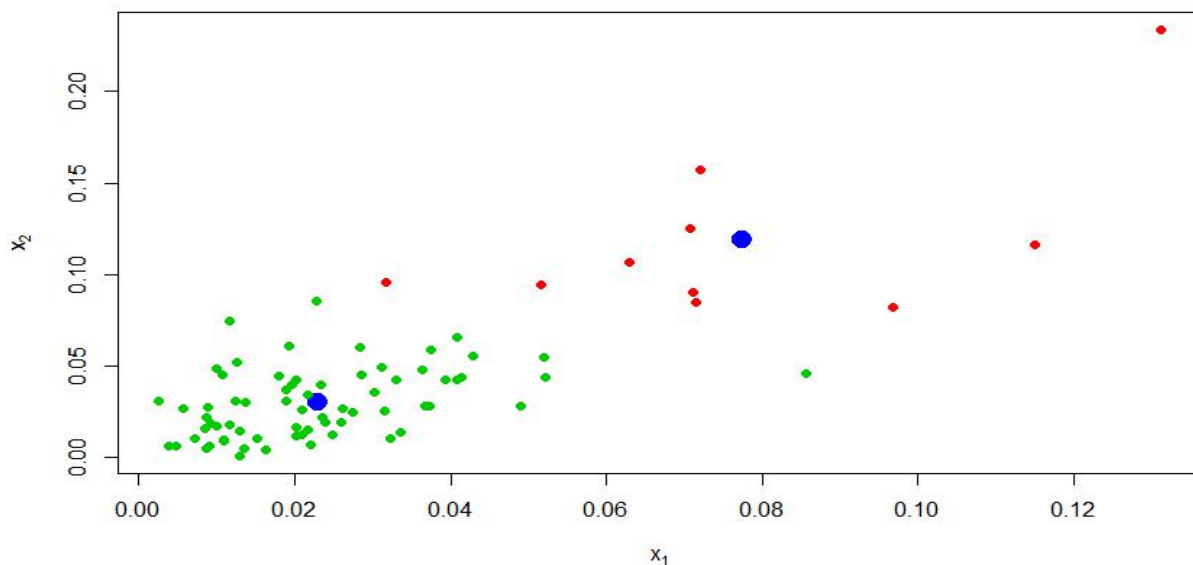


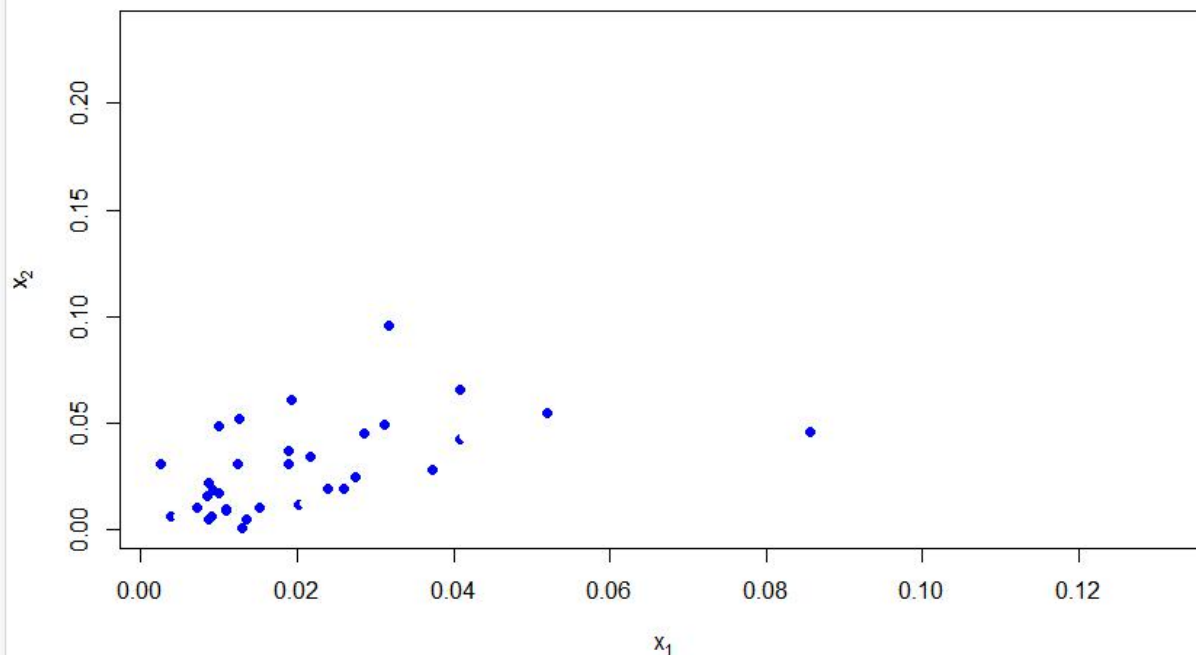
Data Mining Assignment 5

- 1) Read Chapter 8 (Sections 8.1 and 8.2) and Chapter 2 (Section 2.4).
- 2) Use `Kmeans()` with all the default values to find the $k=2$ solution for the first two columns of the sonar test data. Plot these two columns. Also plot the fitted cluster centers using a different color. Finally use the `knn()` function to assign the cluster membership for the points to the nearest cluster center. Color the points according to their cluster membership. Show your R commands for doing so.



- 3) Graphically compare the cluster memberships from the previous problem to the actual labels in the test data. Also compute the misclassification error that would result if you used your clustering rule to classify the data. Show your R commands for doing so.

```
> data <- read.csv("sonar_test.csv", header=FALSE)
>
> x <- data[,1:2]
> plot(x, pch=19, xlab=expression(x[1]), ylab=expression(x[2]))
> y <- data[,61]
> points(x, col=2+2*y, pch=19)
> misscl_err = 1-sum(knnfit==y)/length(y)
> misscl_err
[1] 0.474359
>
```



4) Repeat the previous problem using all 60 columns. Show your R commands for doing so.

```
> data <- read.csv("sonar_test.csv", header=FALSE)
>
> x <- data[,1:60]
> fit <- kmeans(x, 2)
> library(class)
> knnfit <- knn(fit$centers, x, as.factor(c(-1, 1)))
> misscls_err = 1 - sum(knnfit == y) / length(y)
> misscls_err
[1] 0.4358974
>
```

5) Consider the one dimensional data set given $x \leftarrow c(1, 2, 2.5, 3, 3.5, 4, 4.5, 5, 7, 8, 8.5, 9, 9.5, 10)$. Starting with initial cluster center values of 1 and 2 carry out algorithm 10 until convergence by hand for $k=2$ clusters. Show all your work for each step and be sure to say specifically which points are in each cluster at each step.

6) Repeat the previous problem by writing a loop and verify that the final answer is the same and show your R commands for doing so.

```
x <- c(1,2,2.5,3,3.5,4,4.5,5,7,8,8.5,9,9.5,10)
center1 <- 1
center2 <- 2
for (i in 2:10){
  cluster1 <- x[abs(x-center1[i-1])<=abs(x-center2[i-1])]
  cluster2 <- x[abs(x-center1[i-1])>abs(x-center2[i-1])]
  center1[i] <- mean(cluster1)
  center2[i] <- mean(cluster2)
}
```

7) Verify that the kmeans function gives the same solution for the previous problem when you use all of the default values and show your R commands for doing so.

```
> kmeans(x,2)
K-means clustering with 2 clusters of sizes 8, 6

Cluster means:
      [,1]
1 3.187500
2 8.666667

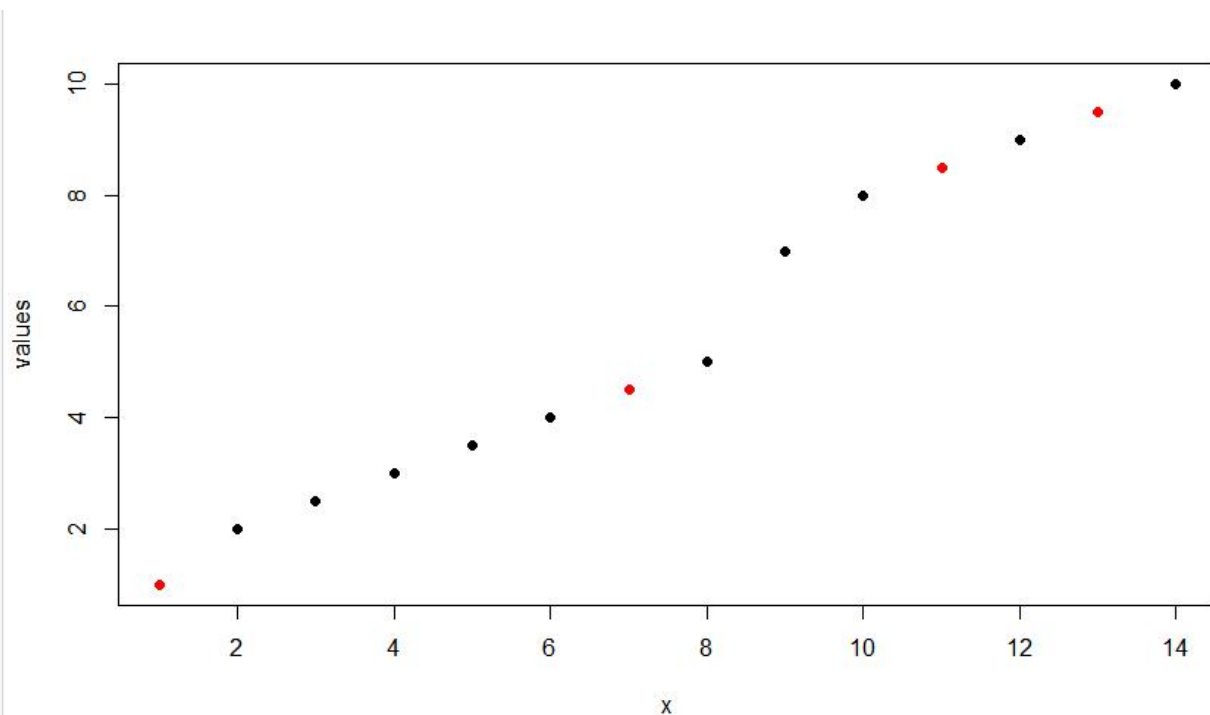
Clustering vector:
[1] 1 1 1 1 1 1 1 1 2 2 2 2 2 2

within cluster sum of squares by cluster:
[1] 12.468750  5.833333
(between_SS / total_SS =  84.9 %)

Available components:

[1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss" "betweenss"    "size"
[8] "iter"         "ifault"

> plot(x, col=fit$cluster, xlab = 'x', ylab='values', pch=19)
```



8) Consider the points $x1 \leftarrow c(1,2)$ and $x2 \leftarrow c(5,10)$.

a) Compute the (Euclidean) distance by hand. Show your work and include a picture of the triangle for the Pythagorean Theorem.

b) Verify that the dist function in R gives the same value as you got in part a. Show your R commands for doing so.

```
> x1 <- c(1,2)
> x2 <- c(5,10)
> data <- (rbind(x1,x2))
> d = dist(data)
> d
      x1
x2 8.944272
> |
```

9) Consider the points $x1 \leftarrow c(1,2,3,6)$ and $x2 \leftarrow c(5,10,4,12)$.

a) Compute the (Euclidean) distance by hand. Show your work.

b) Verify that the dist function in R gives the same value as you got in part a. Show your R commands for doing so.

```
> x1 <- c(1,2,3,6)
> x2 <- c(5,10,4,12)
> data <- (rbind(x1,x2))
> d = dist(data)
> d
      x1
x2 10.81665
> |
```

10) Read Chapter 10.

11) Use a z score cut off of 3 to identify any outliers using the grades for the first midterm at www.stats202.com/spring2008exams.csv. Are there any outliers according to the $z = \pm 3$ rule? What is the value of the largest z score and what is the value of the smallest (most negative) z score? Show your R commands.

```

> data <- read.csv("spring2008exams.csv")
>
> mean_exam <- mean(data[,2],na.rm=TRUE)
> sd_exam <- sd(data[,2],na.rm=TRUE)
> z <- (data[,2]-mean_exam)/sd_exam
> li = sort(z)
> large = li[length(li)]
> small = li[1]
> large
[1] 1.84958
> small
[1] -2.283753
>

```

12) Use a z score cut off of 3 to identify any outliers using the grades for the second midterm at www.stats202.com/spring2008exams.csv. Are there any outliers according to the $z=\pm 3$ rule? What is the value of the largest z score and what is the value of the smallest (most negative) z score? Show your R commands.

```

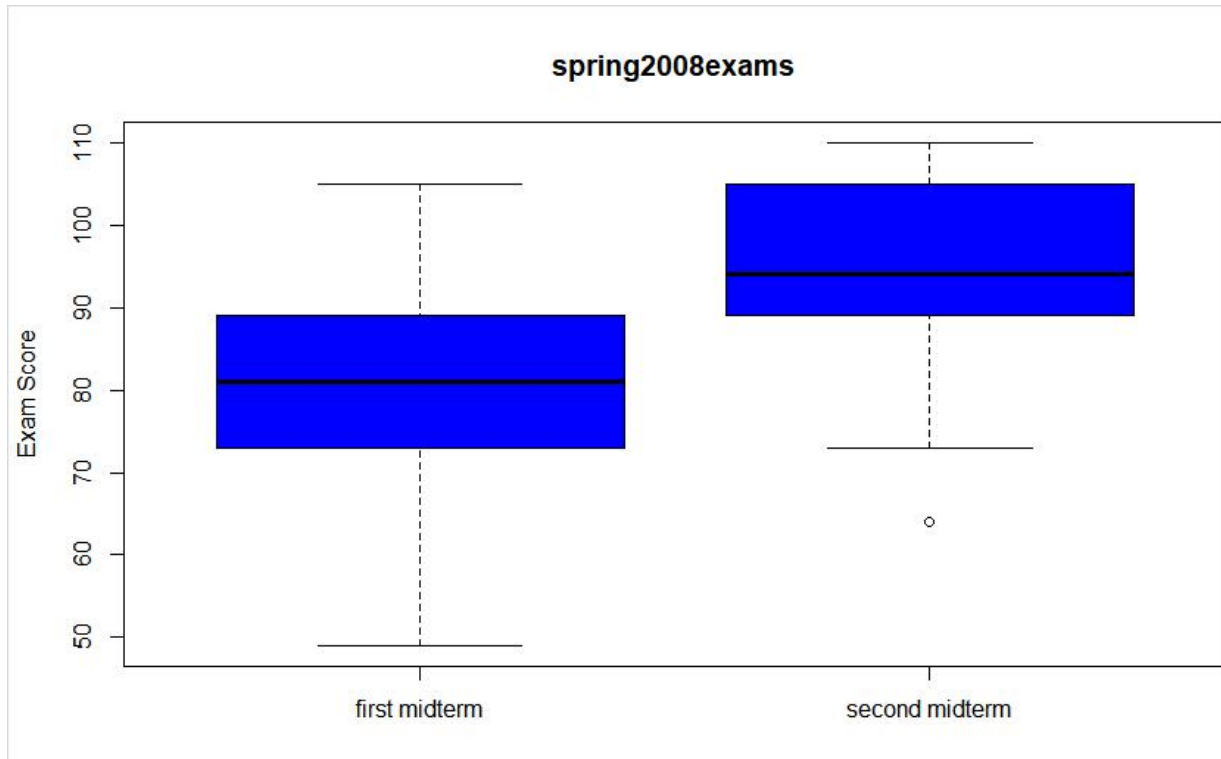
> data <- read.csv("spring2008exams.csv")
>
> mean_exam <- mean(data[,3],na.rm=TRUE)
> sd_exam <- sd(data[,3],na.rm=TRUE)
> z <- (data[,3]-mean_exam)/sd_exam
> li = sort(z)
> large = li[length(li)]
> small = li[1]
> large
[1] 1.299726
> small
[1] -2.396223
>

```

13) Compute the count of each ip address (1st column) in the data stats202log.txt, then use a z score cut off of 3 to identify any outliers for these counts using Excel for the user agent column of the data at www.stats202.com/stats202log.txt. (The user agent column is the second to last column and the value for it in the first row is "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1; .NET CLR 1.1.4322)"). What user agents are identified as outliers using the $z=\pm 3$ rule on the counts of the user agents? What are the z scores for these outliers? (You do not need to show any work for this problem because you are using Excel.)

Ans: Outliner = 0

14) Identify any outliers more than 1.5 IQR's above the 3rd quartile or below the 1st quartile. Verify that these are the same outliers found by the boxplot function using the grades for the second midterm at www.stats202.com/spring2008exams.csv. Show your R commands and include the boxplot. Are any of the grades for the second midterm outliers by this rule? If so, which ones?



15) Use functions to fit a least squares regression model which predicts the exam 2 score as a function of the exam 1 score for the data spring2008exams.csv. Plot the fitted line and determine for which points the fitted exam 2 values are the furthest from the actual values using the model residuals using the midterm grades at www.stats202.com/spring2008exams.csv. Be sure to include the plot. Which student # had the largest POSITIVE residual? Show your R commands.

```
> data<-read.csv("spring2008exams.csv")
>
> model<-lm(data[,3]~data[,2])
> plot(data[,2],data[,3],pch=19,xlab="first midterm",ylab="second midterm",xlim=c(100,200),ylim=c(100,200))
> abline(model)
> max(model$residuals)
[1] 18.17177
```

