

Problem Set 1

QTM 200: Applied Regression Analysis

January 28, 2020

Instructions

- Please show your work! You may lose points by simply writing in the answer. If the problem requires you to execute commands in **R**, please include the code you used to get your answers. Please also include the **.R** file that contains your code. If you are not sure if work needs to be shown for a particular problem, please ask.
- Your homework should be submitted electronically on the course GitHub page in **.pdf** form.
- This problem set is due at the beginning of class on Wednesday, January 22, 2020. No late assignments will be accepted.
- Total available points for this homework is 100.

Question 1 (25 points)

A private school counselor was curious about the average of IQ of the students in her school and took a random sample of 25 students' IQ scores. The following is the data set:

```
1 y <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113, 112, 98,  
      80, 97, 95, 111, 114, 89, 95, 126, 98)
```

Find a 90% confidence interval for the student IQ in the school assuming the population of IQ from which our random sample has been selected is normally distributed.

- ANSWER: The 90% confidence interval is (94.1, 102.7).
- INTERPRETATION: We are 90% confident that that true average IQ of students at her school lies between 94.1 and 102.7.

Question 2 (25 points)

A private school counselor was curious whether the average of IQ of the students in her school is higher than the average IQ score 100 among all the schools in the country. She took a random sample of 25 students' IQ scores. The following is the data set:

```
1 y <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113, 112, 98,
        80, 97, 95, 111, 114, 89, 95, 126, 98)
```

Conduct a test with 0.05 significance level assuming the population of IQ from which our random sample has been selected is normally distributed.

Question 3 (50 points)

Assume y is variable with values 1,2,3,4 standing for “Freshman”, “Sophomore”, “Junior”, and “Senior”, convert y from numbers to characters in R:

```
1 y <- c(1, 2, 1, 3, 4, 1, 1, 4, 2, 1, 3, 4, 3, 2, 1, 3, 4, 1, 2, 3, 1, 1, 2, 1,
        1, 3, 4)
```

Researchers are curious about what affects the education expenditure on public education. The following is available variables in a data set about the education expenditure.

State	50 states in US
Y	per capita expenditure on public education
X1	per capita personal income
X2	Number of residents per thousand under 18 years of age
X3	Number of people per thousand residing in urban areas
Region	1=Northeast, 2= North Central, 3= South, 4=West

Explore the `expenditure` data set and import data into R.

```
1 expenditure <- read.table("expenditure.txt", header=T)
```

- Please plot the relationships among Y , $X1$, $X2$, and $X3$? What are the correlations among them? Describe the graph and the relationships among them.
- Please plot the relationship between Y and $Region$? On average, which region does have the highest per capita expenditure on public education?

- Please plot the relationship between Y and $X1$? Describe this graph and the relationship. Reproduce the above graph including one more variable *Region* and display different regions with different types of symbols and colors.

- R code:

```

1 #####
2 # load libraries
3 # set wd
4 # clear global .envir
5 #####
6
7 # remove objects
8 rm(list=ls())
9 # detach all libraries
10 detachAllPackages <- function() {
11   basic.packages <- c("package:stats", "package:graphics", "package:
      grDevices", "package:utils", "package:datasets", "package:methods", "
      package:base")
12   package.list <- search()[ifelse(unlist(gregexpr("package:", search()))
      ==1, TRUE, FALSE)]
13   package.list <- setdiff(package.list, basic.packages)
14   if (length(package.list)>0) for (package in package.list) detach(
      package, character.only=TRUE)
15 }
16 detachAllPackages()
17
18 # load libraries
19 pkgTest <- function(pkg){
20   new.pkg <- pkg[!(pkg %in% installed.packages()[, "Package"])]
21   if (length(new.pkg))
22     install.packages(new.pkg, dependencies = TRUE)
23   sapply(pkg, require, character.only = TRUE)
24 }
25
26 # here is where you load any necessary packages
27 # ex: stringr
28 # lapply(c("stringr"), pkgTest)
29
30 lapply(c(), pkgTest)
31
32 # set working directory
33 setwd("~/GitHub/QT200Spring2020/problem_sets/PS1")
34
35
36 #####
37 # Problem 1
38 #####
39

```

```

40 y <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113,
      112, 98, 80, 97, 95, 111, 114, 89, 95, 126, 98)
41 z90 <- qnorm((1-.90)/2, lower.tail = FALSE)
42 n <- length(y)
43 mean(y)
44 sample_mean <- mean(y)
45 sample_sd <- sd(y)
46 lower_90 <- sample_mean - (z90 * (sample_sd/sqrt(n)))
47 upper_90 <- sample_mean + (z90 * (sample_sd/sqrt(n)))
48 confint90 <- c(lower_90, upper_90)
49 # ANSWER: 90% confidence interval is (94.1, 102.7)
50
51 #####
52 # Problem 2
53 #####
54
55 y <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113,
      112, 98, 80, 97, 95, 111, 114, 89, 95, 126, 98)
56 #null hypothesis: xbar = 100
57 #alternative hypothesis: xbar > 100
58 #calculate test statistic and p-value
59 mean(y)
60 #mean of sample is 98.44
61 sd(y)
62 #standard deviation of sample is 13.09
63 ts <- ((100-98.44)/(13.09))
64 #sample size
65 n <- length(y)
66 # n=25 therefore df = 24
67 #test statistic = -0.119
68 pt(abs(0.119), df=24, lower.tail=F)
69 # ANSWER: p-value = 0.453 > 0.05 therefore sample mean is not
      significantly greater than 100
70
71 #####
72 # Problem 3
73 #####
74
75 y <- c(1, 2, 1, 3, 4, 1, 1, 4, 2, 1, 3, 4, 3, 2, 1, 3, 4, 1, 2, 3, 1, 1,
      2, 1, 1, 3, 4)
76
77 expenditure <- read.table("expenditure.txt", header=T)
78
79 # PLOT 1
80 plot(expenditure$X1, expenditure$Y,
81       xlab="Per capita personal income", ylab="Per capita expenditure on
      public education")
82 abline(lm(expenditure$Y ~ expenditure$X1))
83 #Moderate positive correlation
84
85 # PLOT 2

```

```

86 plot(expenditure$X2, expenditure$Y,
87       xlab="Number of residents per thousand under 18 years of age", ylab="
      "Per capita expenditure on public education")
88 abline(lm(expenditure$Y ~ expenditure$X2))
89 #Weak negative correlation
90
91 # PLOT 3
92 plot(expenditure$X3, expenditure$Y,
93       xlab="Number of people per thousand residing in urban areas", ylab="
      Per capita expenditure on public education")
94 abline(lm(expenditure$Y ~ expenditure$X3))
95 #Weak positive correlation
96
97
98 # PLOT 4
99
100 expenditure$fourregions <- NA
101 expenditure$fourregions <- factor(NA, levels=c("Northeast", "North
      Central", "South", "West"))
102 expenditure$fourregions[expenditure$Region==1] <- "Northeast"
103 expenditure$fourregions[expenditure$Region==2] <- "North Central"
104 expenditure$fourregions[expenditure$Region==3] <- "South"
105 expenditure$fourregions[expenditure$Region==4] <- "West"
106 boxplot(expenditure$Y ~ expenditure$fourregions,
107          xlab= "Region",
108          ylab= "Per capita expenditure on public education")
109
110 # On average, the Western region has the highest per capita expenditure
      on public education
111
112 # PLOT 5
113
114 plot(expenditure$X1, expenditure$Y,
115       xlab="Per capita personal income", ylab="Per capita expenditure on
      public education")
116 abline(lm(expenditure$Y ~ expenditure$X1))

```