

Problem Set 1

QTM 200: Applied Regression Analysis

Due: January 29, 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.

```
1 #find t statistics for a 90% confidence level  
2 t90 <- qt((1-0.9)/2, df= 24, lower.tail = FALSE)  
3 n <- length(y) #sample size  
4 ymean <- mean(y) # sample mean  
5 ysd <- sd(y) # sample standard deviation  
6 lower <- ymean - t90*(ysd/sqrt(n))
```

```

7 upper <- ymean + t90*(ysd/sqrt(n))
8 #construct the confidence interval
9 confint90 <- c(lower, upper)
10 confint90 #[1] 93.95993 102.92007

```

We are 90% confident that the population IQ lies within the range of 93.95993 to 102.92007

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.

```

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)
2 #T test assuming that the school average IQ is higher than the country average
  (100)
3 t.test(y, mu=100, alternative="greater", conf.level=0.95)
4 #t = -0.59574, df = 24, p-value = 0.7215

```

Since p-value = 0.7215, which is greater than significance level 0.05, This does not reject the null hypothesis and thus we cannot conclude that the school mean IQ is higher than the country's mean IQ

Question 3 (50 points)

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 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)
```

- Please plot the relationships among Y , $X1$, $X2$, and $X3$? What are the correlations among them (you just need to describe the graph and the relationships among them)?

As we can see in *Figure 1* and *Figure 3*, public education expenditure is positively correlated with both income and urban population; As seen in *Figure 2*, education expenditure is weakly negatively correlated with number of residents at school age.

- Please plot the relationship between Y and *Region*? On average, which region has the highest per capita expenditure on public education?

As seen in *Figure 4*, All four regions have very different expenditures on public education. The West has by far the highest per capita expenditure on public education and the three other regions are all comparatively lower.

- 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.

As seen in *Figure 5*, the public educational expenditure is positively correlated to personal income per capita. Generally, there are no evident clusters for any of the regions, but differences in regional distribution are still distinct. We can also see from the distribution of the green plus sign that the South contributed to this correlation the most, with other regions spreading more to the upper right.

```
1 #Relationship between Y,X1,X2 and X3
2 plot(expenditure$X1, expenditure$Y, main = "Public education Expenditure
        vs. personal income",
3                                     xlab = "personal income per capita",
4                                     ylab = "expenditure per capita")
5 #The best fit line
6 abline(lm(expenditure$Y~expenditure$X1))
7
8 #X2
9 plot(expenditure$X2, expenditure$Y, main = "Public education Expenditure
        vs. underaged population",
10                                     xlab = "number of resident per
        thousands under 18",
11                                     ylab = "expenditure per capita")
12 #The best fit line
```

```

13 abline(lm(expenditure$Y~expenditure$X2))
14 #X3
15 plot(expenditure$X3, expenditure$Y, main = "Public education Expenditure
    vs. urban population",
16       xlab = "number of urban resident per
    thousands",
17       ylab = "expenditure per capita")
18 #The best fit line
19 abline(lm(expenditure$Y~expenditure$X3))
20
21 #Relationship between Y and Regions
22 library(Rmisc)
23 summarySE(data = expenditure, measurevar = "Y", groupvars = "Region")
24 #Create a new factor variable YR
25 expenditure$YR <- expenditure$Y
26 expenditure$YR <- factor(NA, levels = c("Northeast", "North Central", "
    South", "West"))
27 expenditure$YR[expenditure$Region == 1] <- "Northeast"
28 expenditure$YR[expenditure$Region == 2] <- "North Central"
29 expenditure$YR[expenditure$Region == 3] <- "South"
30 expenditure$YR[expenditure$Region == 4] <- "West"
31 summary(expenditure$YR)
32 #A side by side boxplot
33 boxplot(expenditure$Y~expenditure$YR, main = "Public education
    Expenditure vs. Regions",
34         xlab = "Regions",
35         ylab = "expenditure per capita")
36 #Y vs X1 plus regional information
37 plot(expenditure$X1, expenditure$Y, main = "Public education Expenditure
    vs. personal income",
38       xlab = "personal income per capita",
39       ylab = "expenditure per capita", col = expenditure$Region, pch =
    expenditure$Region)
40 legend(x="topright", legend = levels(expenditure$YR), col = c(1,2,3,4),
    pch= c(1,2,3,4))

```

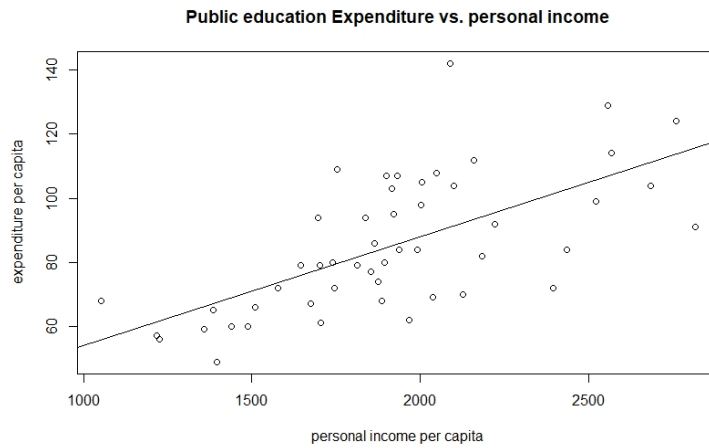


Figure 1: Y and X1

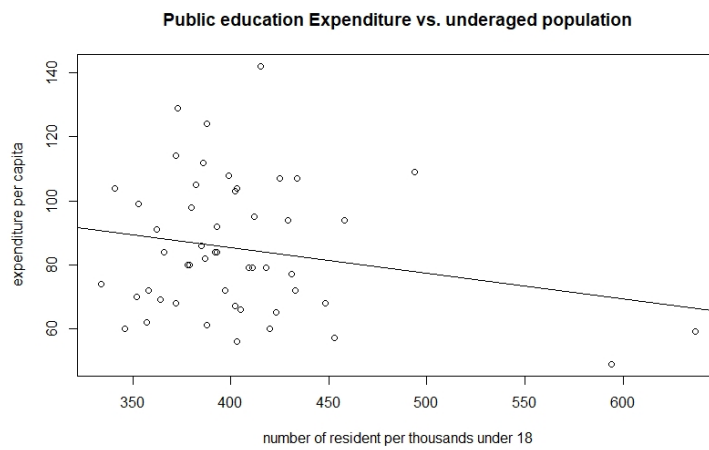


Figure 2: Y and X2

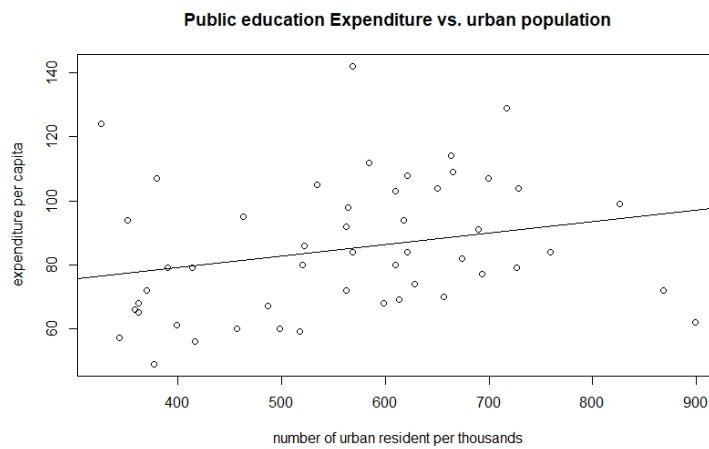


Figure 3: Y and X3

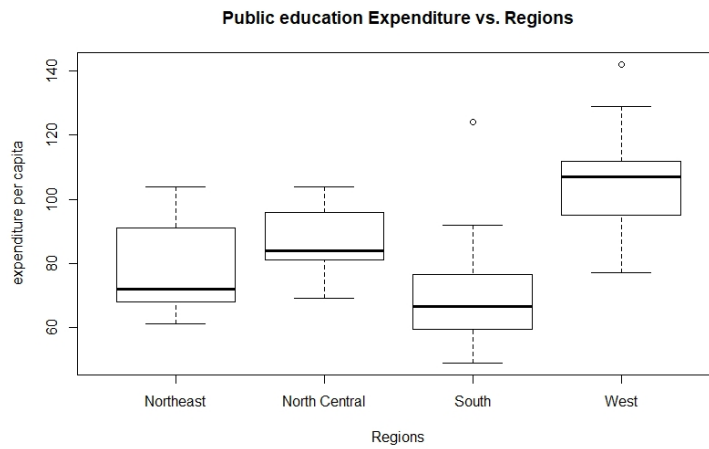


Figure 4: Y and Region

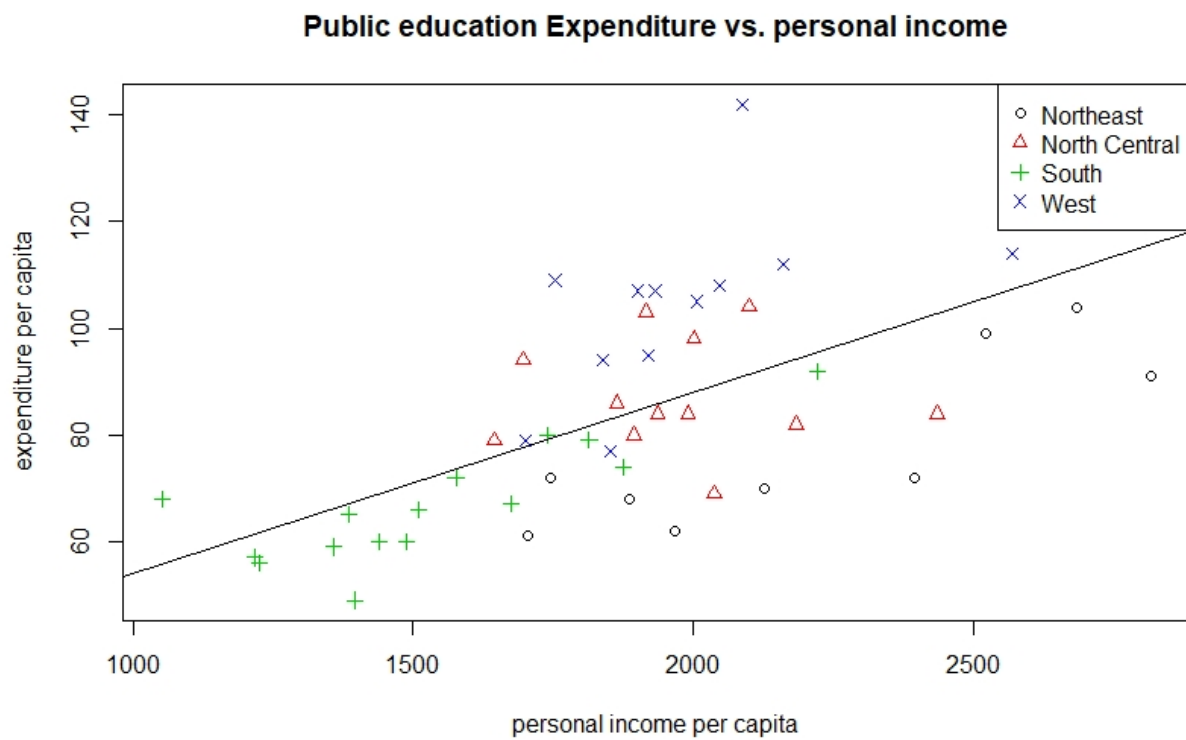


Figure 5: Y vs. X1 with Regions labeled