

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:

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 #####  
2 # Problem 1  
3 #####  
4  
5 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)  
6 #Find the sum  
7 sum(y)  
8 #Find the mean
```

```

9 sum(y)/length(y)
10 sample_mean=mean(y)
11 #Find the sum of errors
12 demeanedSumSimple = y - mean(y)
13 demeanedSumSimple
14 #Find the squared error
15 squaredError= demeanedSumSimple^2
16 squaredError
17 #Find the variance
18 variance=sum(squaredError)/length(y)
19 variance
20 #Find the standard deviation
21 Sd= sd(y,na.rm=FALSE)
22 sample_sd=Sd
23 #Given confidence coefficient = 0.9
24 z90=qt((1-0.9)/2,df=length(y)-1, lower.tail = FALSE)
25 n= length(y)
26 lower_90 = sample_mean-(z90*(sample_sd/sqrt(n)))
27 upper_90 = sample_mean+(z90*(sample_sd/sqrt(n)))
28 confint90 = c(lower_90,upper_90)
29 # [93.95993 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
2 #####
3 # Problem 2
4 #####
5
6 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)
7
8 # set Ho:mu < 100
9 # set Ha:mu >= 100
10
11 mu = 100
12 #Find the sum

```

```

13 sum(y)
14 #Find the mean
15 mean(y)
16 #Find the sum of errors
17 demeanedSumSimple = y - mean(y)
18 demeanedSumSimple
19 sum(demeanedSumSimple)
20 #Find the squared error
21 squaredError= demeanedSumSimple^2
22 #Find the variance
23 variance=sum(squaredError)/length(y)-1
24 variance
25
26 z=(mean(y)-mu)/variance
27
28 critical_value = qt(0.05, df=length(y)-1, lower.tail = T )
29
30 if(critical_value<z){
31   print("there is not enough evidence to disprove Ho, the null hypothesis,
32     which states that students in her school is lower than the average IQ
33     score 100 among all the schools in the country")
34 }
35 #There is not enough evidence to disprove Ho, the null hypothesis, which
36   states that students in her school is lower than the average IQ score 100
37   among all the schools in the country, thus the hypothesis testing suggest
38   the counselor's hypothesis might not be right

```

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

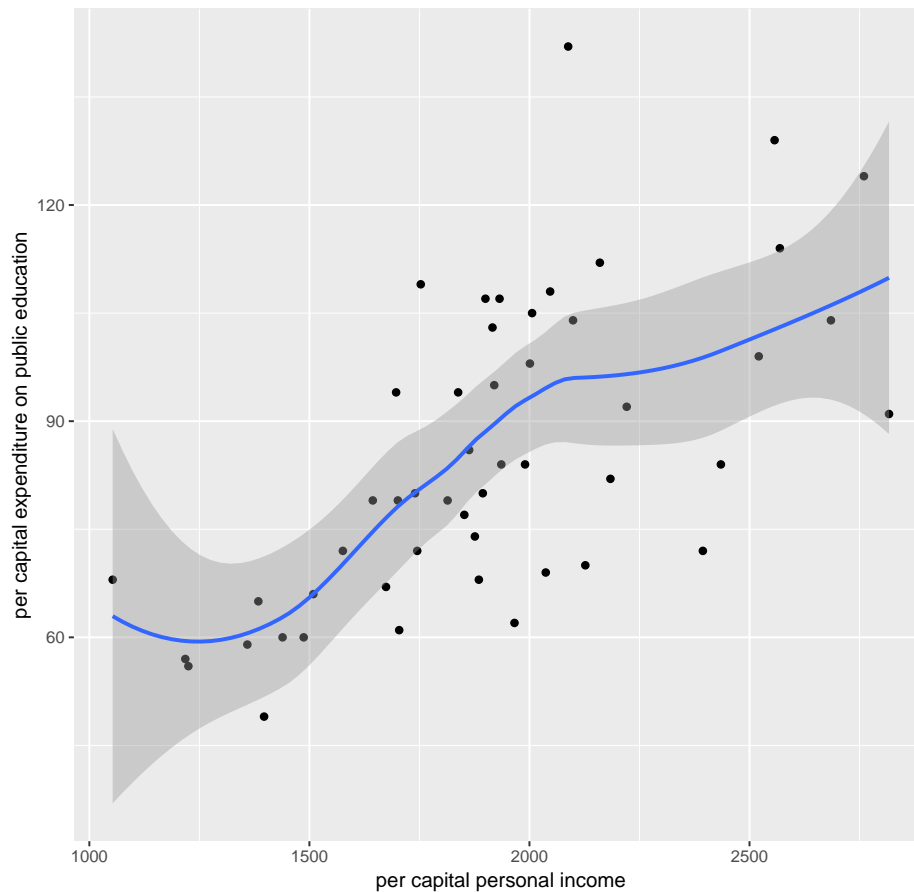
Explore the `expenditure` data set and import data into R.

```

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

```

Figure 1: Relationship between per capital personal income and per capita expenditure on public education



- 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)?

```

1
2 #1
3
4 #Y & X1
5 pdf("plot1.pdf")
6 ggplot(expenditure, aes(x=X1, y=Y)) + geom_point() + labs(x="per capital
  personal income", y="per capital expenditure on public education",
  title="")+geom_smooth()
7 dev.off()
8 #Observation: The relationship between X1 and Y is almost linear and the
  line is upwards sloping. With per capital personal income increases,

```

Figure 2: Relationship between number of residents per thousand under 18 years of age and per capital expenditure on public education

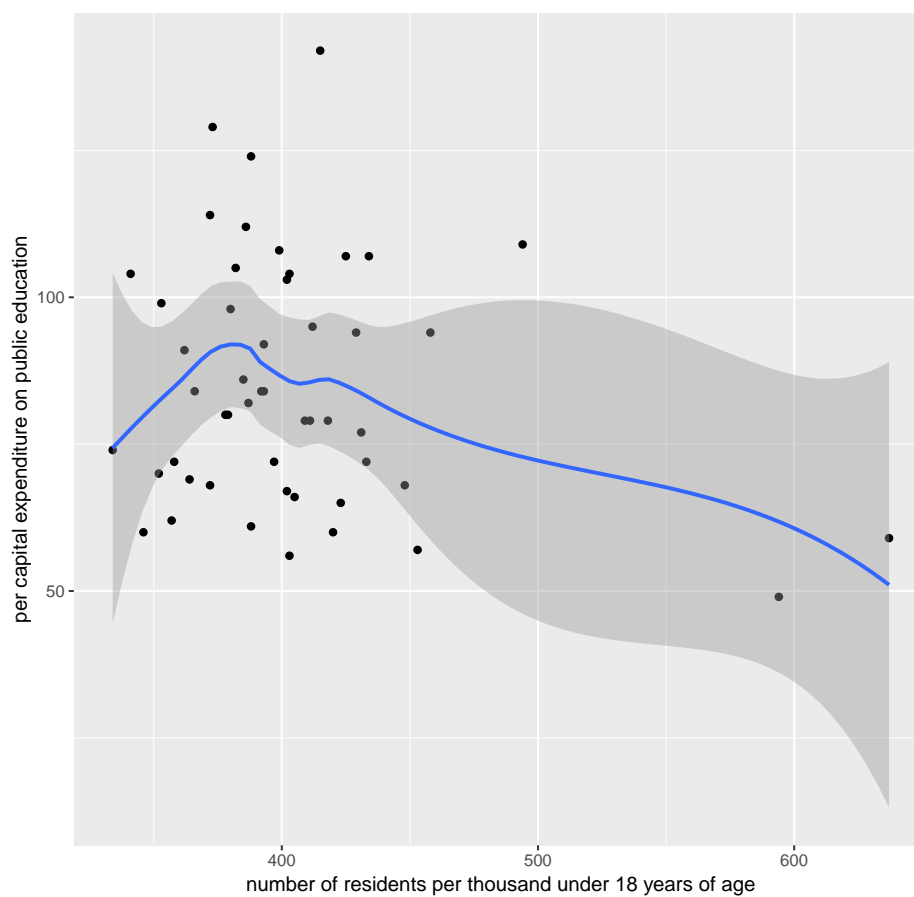
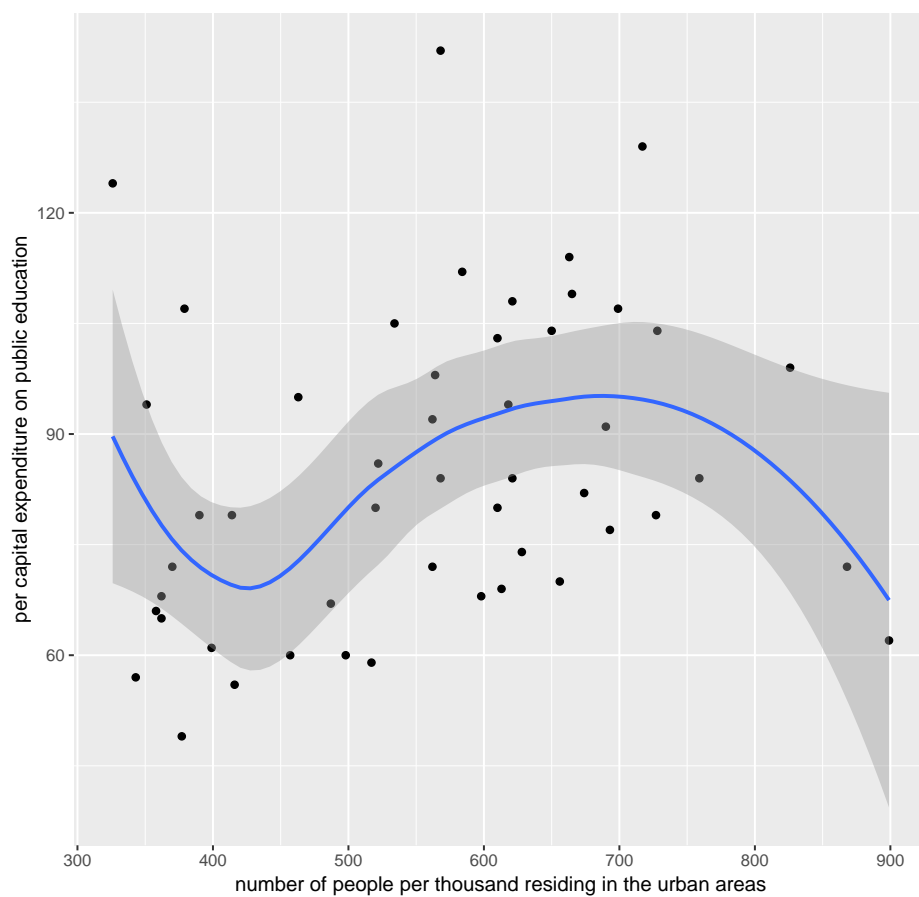


Figure 3: Relationship between number of people per thousand residing in the urban areas and per capital expenditure on public education



```

9     per capital expenditure on public education also increases.
10 #Y & X2
11 pdf("plot2.pdf")
12 ggplot(expenditure, aes(x=X2, y=Y)) + geom_point() + labs(x="number of
    residents per thousand under 18 years of age", y="per capital
    expenditure on public education", title="")+geom_smooth()
13 dev.off()
14 #Observation: The relationship between X2 and Y is almost linear and the
    line is downwards sloping. With number of residents per thousand
    under 18 years of age increases, per capital expenditure on public
    education also decreases.
15
16 #Y & X3
17 pdf("plot3.pdf")
18 ggplot(expenditure, aes(x=X3, y=Y)) + geom_point() + labs(x="number of
    people per thousand residing in the urban areas", y="per capital
    expenditure on public education", title="")+geom_smooth()
19 dev.off()
20 #Observation: The relationship between X3 and Y is not linear because the
    curve has 3 inflection points. When X3 approaches to 400, Y value is
    decreasing; when X3 approaches to 700, Y value is increasing; when X3
    increases after 700, Y value decreases.

```

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

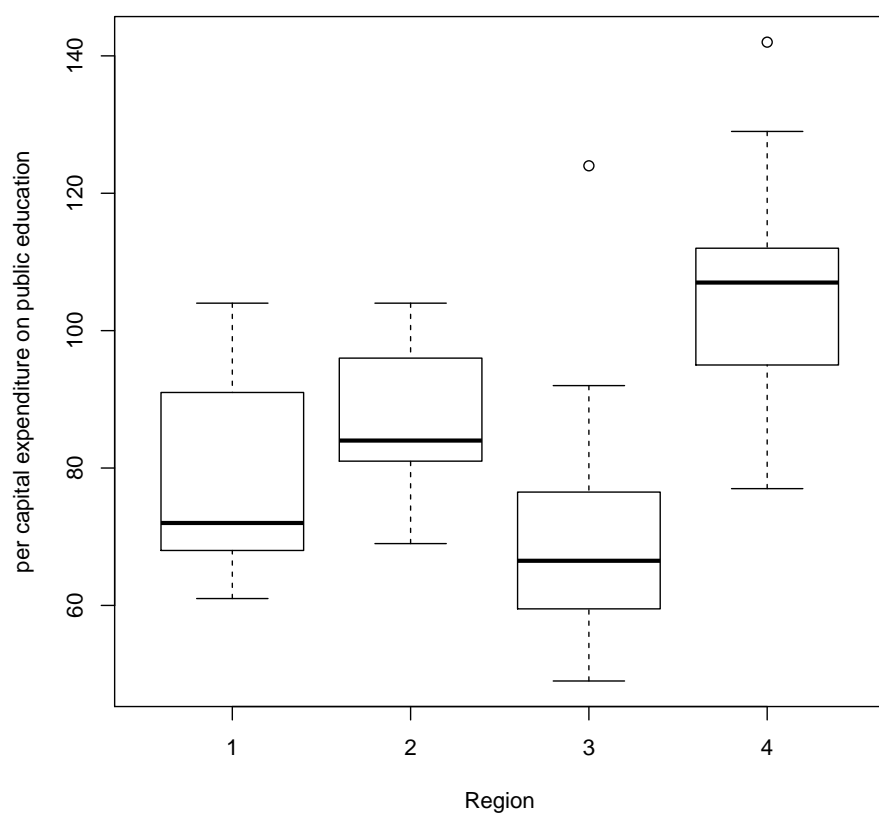
```

1
2 #2(graphics for this question is under #3 for no reason: I can't fix it
   :( )
3 pdf("plot4.pdf")
4
5 input=expenditure[,c('Region', 'Y')]
6 print(head(input))
7 boxplot(Y~Region, data = expenditure, xlab = "Region",
8 ylab = "per capital expenditure on public education", main = "")
9
10 dev.off()
11
12 #Observation: I use the box plot to compare across regions in terms of
    per capital expenditure on public education. According to the chart
    we find on average West(4) region has the largest per capital
    expenditure on public education whereas South(3) region has the
    lowest

```

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

Figure 4: Relationship between Regions and per capital expenditure on public education



different regions with different types of symbols and colors.

```
1
2 #3
3
4 pdf("plot5.pdf")
5
6 ggplot(expenditure, aes(x=X1, y=Y))+ geom_point()+labs(x="per capital
  personal income", y="per capital expenditure on public education",
  title="")+ geom_smooth()
7 #Observation: The relationship between X1 and Y is almost linear and the
  line is upwards sloping. With per capital personal income increases,
  per capital expenditure on public education also increases.
8 dev.off()
9
10 pdf("plot6.pdf")
11 ggplot(expenditure, aes(x=X1, y=Y))+geom_point(aes(color=Region))+ labs(x
  ="per capital personal income", y="per capital expenditure on public
  education", title="")
12 dev.off()
```

Figure 5: Relationship between per capital personal income and per capital expenditure on public education

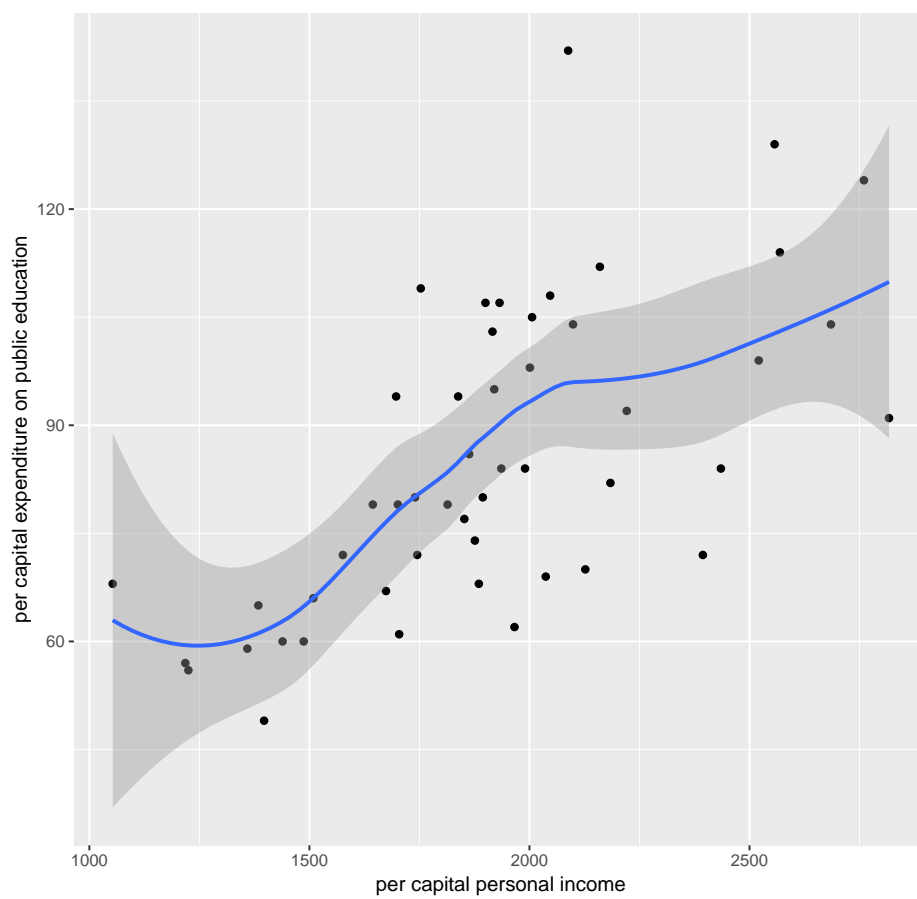


Figure 6: Relationship between per capital personal income and per capital expenditure on public education with Region on the side

