

# Problem Set 1

## QTM 200: Applied Regression Analysis

Due: January 27, 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 ##finding the confidence interval  
2 #first the mean was found  
3 m <- mean(y)  
4 #then the standard deviation  
5 s <- sd(y)  
6 #n is the number of samples
```

```

7 n <- 25
8 #finding the error using t distribution as n is less than 30
9 error <- qt(0.95, df=n-1)*s/sqrt(n)
10 #left confidence interval
11 left <- m-error
12 #right confidence interval
13 right <- m+error
14 left
15 right

```

In order to find the confidence interval, the mean and standard deviation of the data was determined and stated as "m" and "s" respectively. The mean came out to be 98.44 and the standard deviation came out to be 13.09. Then I assigned "n" as the total number of sample, which is 25. With these variables, the error could be determined by the use of a t-distribution as the number of samples is less than 30, thus illustrating a small sample. With the confidence interval being 90%, we determined the 5% error on each side of the t-distribution. The error came out to be 4.48. In order to determine the lower confidence interval, the difference between mean and error was determined. To determine the upper confidence interval, the sum of mean and error was determined.

It was found that the 90% confidence interval for the student IQ in the school is 93.96 and 102.92

## 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 #null hypothesis Ho: pi < pio
2 pio <- 100
3 pi <- mean(y)
4 pi
5 ts <- qt(.95, df=n-1)*s/sqrt(n)
6 ts
7 p = pt(abs(ts), df=n-1, lower.tail=F)
8 p

```

To conduct a test with 0.05 significance level, first the assumptions need to be stated. The data is continuous, random sampling and the sample is distributed normally.

The null hypothesis is  $\mu < 100$  and the alternative hypothesis is  $\mu \geq 100$ .

The population mean is 100, while the sample mean is 98.44. The test statistic is determined by using a t-distribution with the confidence interval being 0.95 as this is a one-sided test. The t-distribution is used as this is a small sample (less than 30). The test statistic came out to be 4.48 (which is identical to the error of Question 1) and with this, the p-value is determined to be  $7.79 \times 10^{-5}$ , which is smaller than 0.05, and thus leading to rejecting the null hypothesis. We can state that the average IQ of the students in her school is equal or higher to the average IQ score among the population.

## 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.

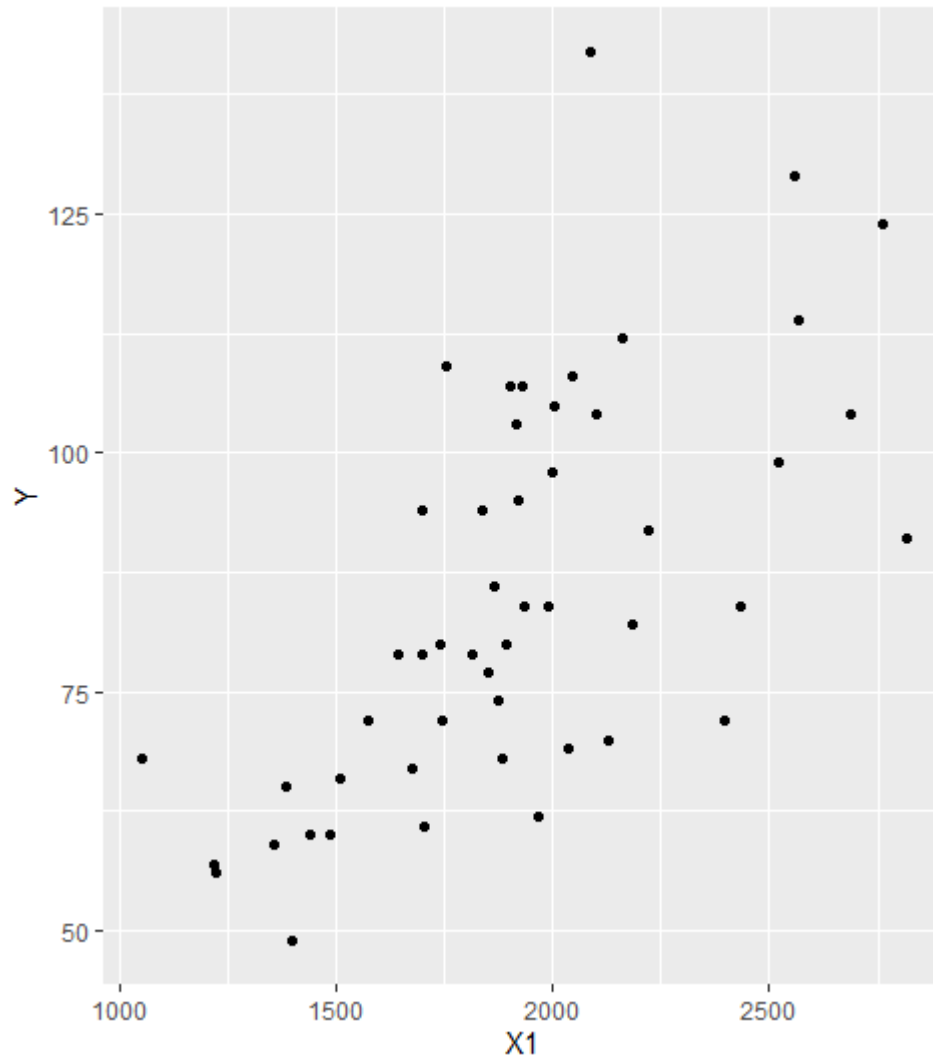
State	50 states in US
Y	per capita expenditure on public education
X1	per capita personal income
X2	Number of residents per thousand under 18 years of age
X3	Number of people per thousand residing in urban areas
Region	1=Northeast, 2= North Central, 3= South, 4=West

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

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

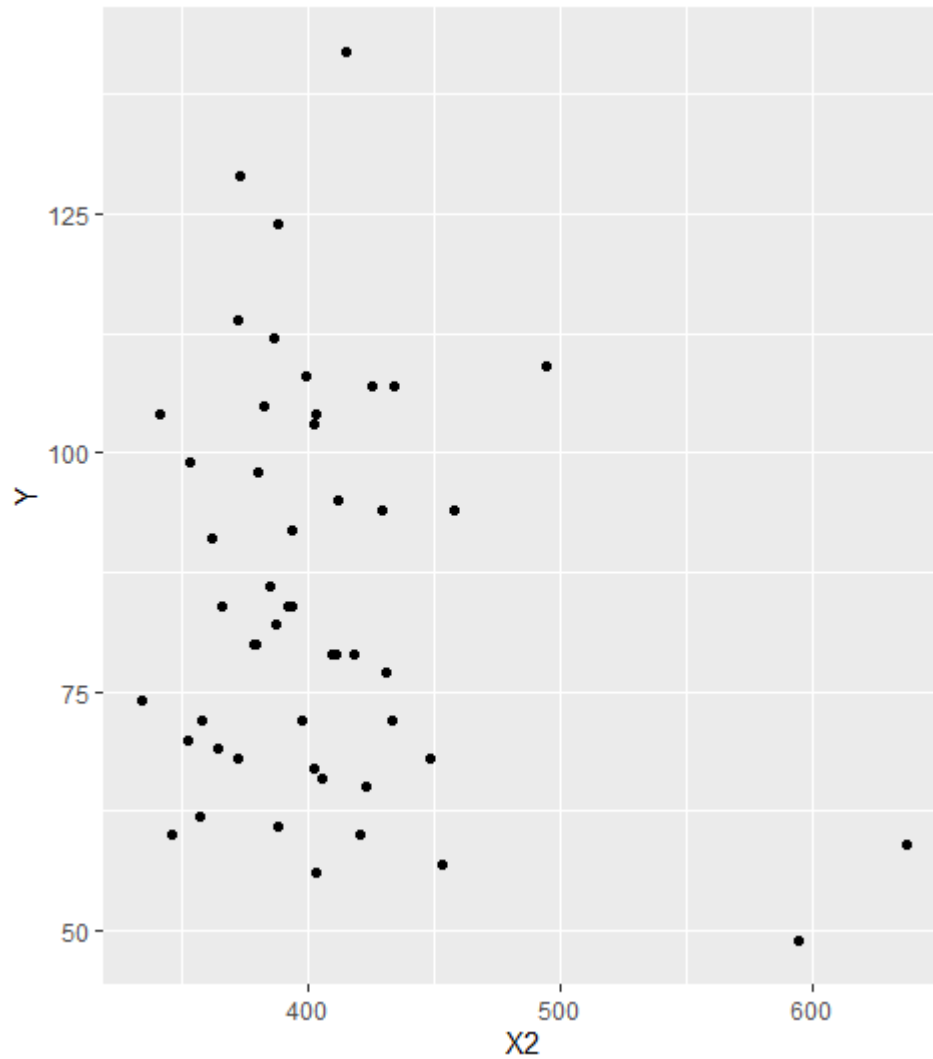
- Please plot the relationships among  $Y$ ,  $X1$ ,  $X2$ , and  $X3$ ? What are the correlations among them? Describe the graph and the relationships among them.

```
1 ggplot(expenditure, aes(x=X1, y=Y)) +
2   geom_point()
```



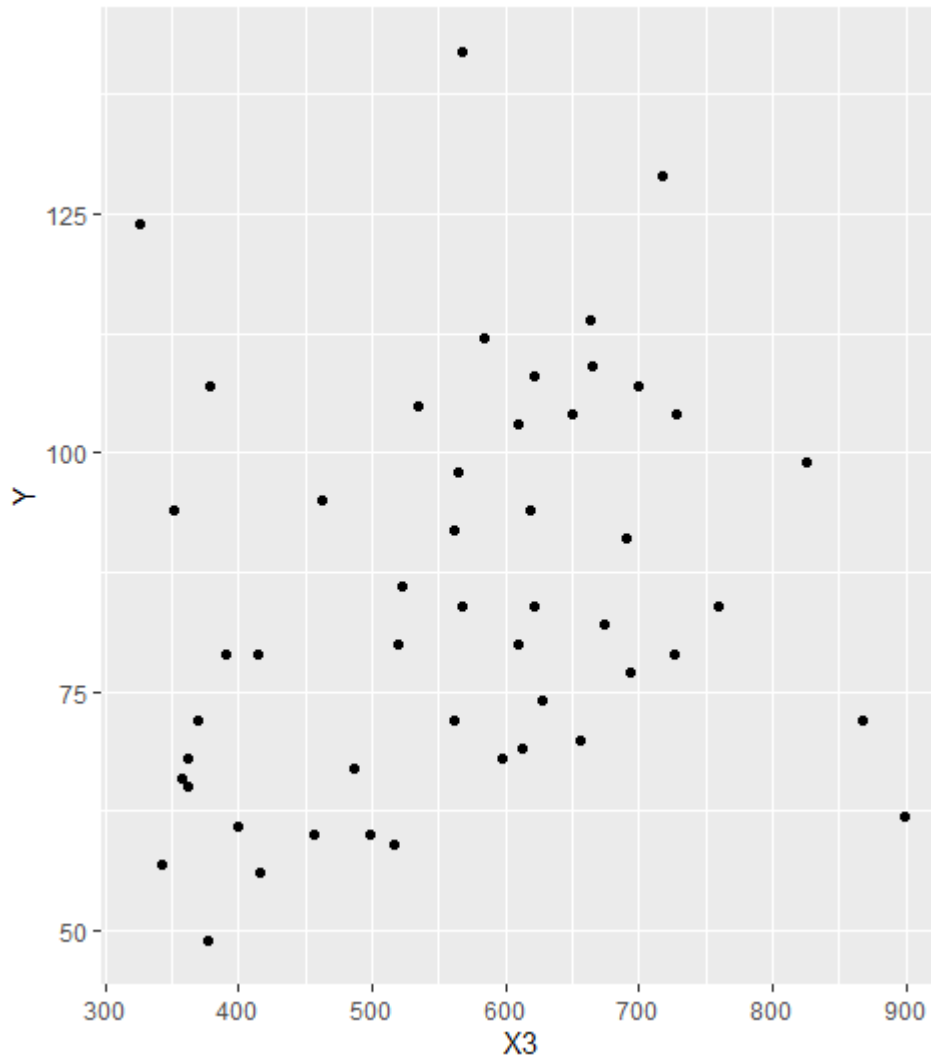
This graph illustrates a slight positive correlation showing that as per capita personal income increases, the capita expenditure on public education also increases.

```
1 ggplot(expenditure , aes(x=X2, y=Y))+  
2   geom_point()
```



This graph illustrates no correlation, as most of the data is on the lower end of the x axis but the Y axis is spread. This could be said that the population with more minors by proportion would end up actually spending expenditure on public education as there is a population catered towards them. The more residents per thousand under 18 years of age, the less or will mostly likely not spend on public education.

```
1 ggplot(expenditure , aes(x=X3, y=Y)) +
2   geom_point()
```

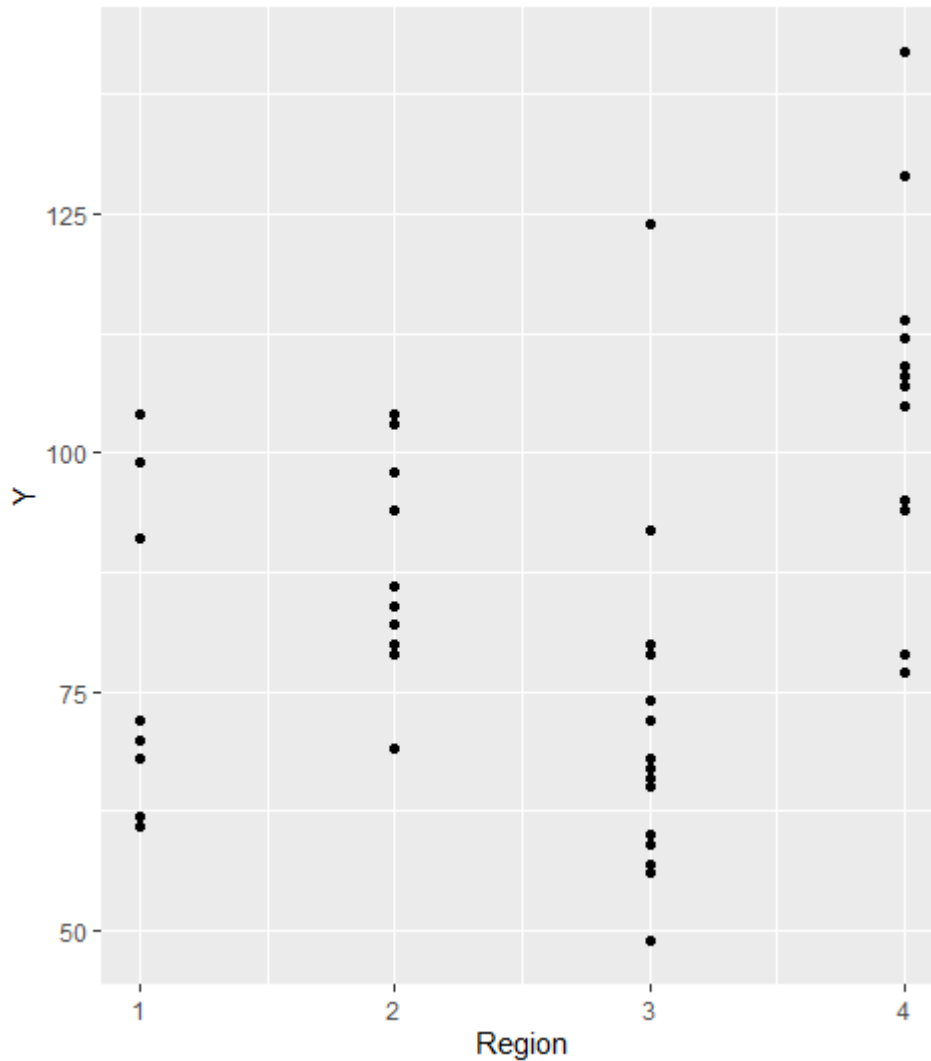


This graph is shown to also have no correlation, however the data is more spread out throughout the y and x axis. This illustrates that there is no correlation between the number of people per thousand residing in urban areas to the expenditure on public education.

The three graphs illustrate that the higher personal income within the region, the more expenditure on public education. The regions with the most minors tend to have expenditure on public education, while the regions with less minors have shown to have low or no expenditure, all regardless of whether the region is urban or not.

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

```
1 ggplot(expenditure, aes(x=Region, y=Y)) +
2   geom_point()
```



Region 4 on average has the highest per capital expenditure on public education compared to the other regions. Not only is the min and max the largest in Region 4, but on average, Region 4 is a lot higher per capital expenditure on public education compared to other regions.

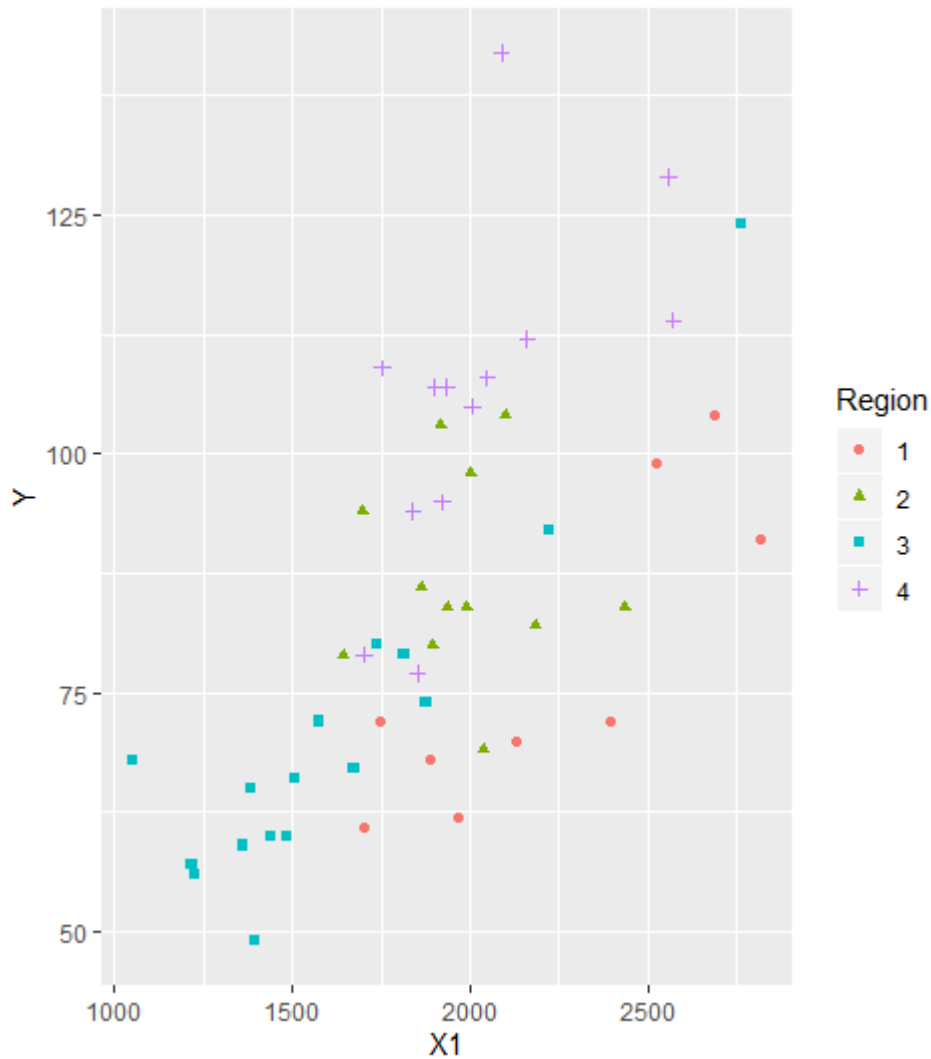
- Please plot the relationship between  $Y$  and  $X1$ ? Describe this graph and the relationship. Reproduce the above graph including one more variable *Region* and display different regions with different types of symbols and colors.

```

1 expenditure[, 'Region'] <- factor(expenditure[, 'Region'])
2 class(expenditure$Region)
3
4 ggplot(expenditure, aes(x=X1, y=Y, color=Region, group=Region))+
5   geom_point(aes(shape=Region))

```

I had to change the Region from a continuous variable to a factor as it could not categorize Region into shape due to the fact that the data type was continuous.



The graph illustrate that in X2 data set, the Region 4 seems to have the highest per capita expenditure on public education in addition to having a higher capita of personal income. It is shown that Region 3 seems to have the lowest capita of personal income on average, and thus also have the lowest average in expenditure on public education.