

# 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.

#### Solution

R code:

```
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 z90 <- qnorm((1 - .90)/2, lower.tail = FALSE)
```

```

3 n <- length(y)
4 sample_mean <- mean(y)
5 sample_sd <- sd(y)
6 lower_90 <- sample_mean - (z90 * (sample_sd/sqrt(n)))
7 upper_90 <- sample_mean + (z90 * (sample_sd/sqrt(n)))
8 confint90 <- c(lower_90, upper_90)
9 confint90

```

Results: [94.13283, 102.74717]

## 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.

### Solution

R code:

```

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(y, mu = 100,
3       alternative = "greater")
4 # OR
5 mu <- 100
6 n <- length(y)
7 sample_mean <- mean(y)
8 sample_sd <- sd(y)
9 teststatistic <- (sample_mean-mu)/(sample_sd/sqrt(n))
10 teststatistic
11 pt(abs(teststatistic), df = n - 1)

```

### Results

Due to the p-value of 0.7215 being larger than the required significance level of 0.05, we fail to reject the null hypothesis and cannot accept the alternative hypothesis. The evidence does not supports the hypothesis that the average IQ students in the private school is higher than the average IQ score 100 among all the schools in the country.

## 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.

### Solution to converting y from numbers to characters

R code:

```
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)
2 y.characters <- c("Freshman", "Sophomore", "Junior", "Senior")[y]
3 y.characters
```

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 (you just need to describe the graph and the relationships among them)?

R code:

```
1 expenditure <- read.table("expenditure.txt", header=T)
2 head(expenditure,6)
3 cor(expenditure$Y, expenditure$X1)
4 exp_data <- expenditure[,2:length(expenditure)]
5 round(cor(exp_data),2)
6
7 plot(expenditure$X1, expenditure$Y)
8 plot(expenditure$X2, expenditure$Y)
9 plot(expenditure$X3, expenditure$Y)
10
11 plot(expenditure$X1, expenditure$X2)
12 plot(expenditure$X1, expenditure$X3)
13
14 plot(expenditure$X2, expenditure$X3)
```

After plotting the relationships and building a correlation matrix among  $Y$ ,  $X1$ ,  $X2$ , and  $X3$ . I've come up with these observations.

- |                           |             |  |                             |
|---------------------------|-------------|--|-----------------------------|
|                           | From        |  | <i>linear association</i>   |
|                           | Direction   |  | <i>positive association</i> |
| • Between $Y$ and $X1$ :  | Strength    |  | <i>moderately strong</i>    |
|                           | Outliers    |  | <i>present</i>              |
|                           | Correlation |  | <i>0.65</i>                 |
|                           |             |  |                             |
|                           | From        |  | <i>linear association</i>   |
|                           | Direction   |  | <i>negative association</i> |
| • Between $Y$ and $X2$ :  | Strength    |  | <i>weak</i>                 |
|                           | Outliers    |  | <i>present</i>              |
|                           | Correlation |  | <i>-0.21</i>                |
|                           |             |  |                             |
|                           | From        |  | <i>linear association</i>   |
|                           | Direction   |  | <i>positive association</i> |
| • Between $Y$ and $X3$ :  | Strength    |  | <i>weak</i>                 |
|                           | Outliers    |  | <i>present</i>              |
|                           | Correlation |  | <i>0.25</i>                 |
|                           |             |  |                             |
|                           | From        |  | <i>linear association</i>   |
|                           | Direction   |  | <i>negative association</i> |
| • Between $X1$ and $X2$ : | Strength    |  | <i>moderately strong</i>    |
|                           | Outliers    |  | <i>present</i>              |
|                           | Correlation |  | <i>-0.53</i>                |
|                           |             |  |                             |
|                           | From        |  | <i>linear association</i>   |
|                           | Direction   |  | <i>positive association</i> |
| • Between $X1$ and $X3$ : | Strength    |  | <i>moderately strong</i>    |
|                           | Outliers    |  | <i>present</i>              |
|                           | Correlation |  | <i>0.60</i>                 |
|                           |             |  |                             |
|                           | From        |  | <i>linear association</i>   |
|                           | Direction   |  | <i>negative association</i> |
| • Between $X2$ and $X3$ : | Strength    |  | <i>weak</i>                 |
|                           | Outliers    |  | <i>present</i>              |
|                           | Correlation |  | <i>-0.37</i>                |
- Please plot the relationship between  $Y$  and *Region*? On average, which region has the highest per capita expenditure on public education?

```

1 plot(expenditure$Region, expenditure$Y)
2
3 region1 <- filter(expenditure, expenditure$Region == 1)
4 mean(region1$Y)
5
6 region2 <- filter(expenditure, expenditure$Region == 2)

```

```

7 mean(region2$Y)
8
9 region3 <- filter(expenditure, expenditure$Region == 3)
10 mean(region3$Y)
11
12 region4 <- filter(expenditure, expenditure$Region == 4)
13 mean(region4$Y)

```

On average, Region 4 has 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.

```

1 plot(expenditure$Region, expenditure$Y)
2
3 region1 <- filter(expenditure, expenditure$Region == 1)
4 mean(region1$Y)
5
6 region2 <- filter(expenditure, expenditure$Region == 2)
7 mean(region2$Y)
8
9 region3 <- filter(expenditure, expenditure$Region == 3)
10 mean(region3$Y)
11
12 region4 <- filter(expenditure, expenditure$Region == 4)
13 mean(region4$Y)

```

Between $Y$ and $X1$ :	From	<i>linear association</i>
	Direction	<i>positive association</i>
	Strength	<i>moderately strong</i>
	Outliers	<i>present</i>
	Correlation	<i>0.65</i>

