

Problem Set 1

Applied Stats/Quant Methods 1

Due: October 1, 2021

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 GitHub in .pdf form.
- This problem set is due before 8:00 on Friday October 1, 2021. No late assignments will be accepted.
- Total available points for this homework is 100.

Question 1 (50 points): Education

A 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 # Problem 1
2 #####
3
4 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)
```

1. Find a 90% confidence interval for the average student IQ in the school.

Answer:

```

1 n <- 25    ## sample size
2 SAMMEAN <- mean(y)    ## calculate sample mean
3 SAMSD <- sd(y)    ## calculate sample standard deviation
4 MARGIN <- qt(0.95, df=n-1)*SAMSD/sqrt(n)
5 LL <- SAMMEAN - MARGIN ## Lower Limit
6 LL
7
8 HL <- SAMMEAN + MARGIN ## Upper Limit
9 HL

```

Confidence interval for the average student IQ: lower Limit= 93.96, Upper Limit = 102.92

2. Next, the school counselor was curious whether the average student IQ in her school is higher than the average IQ score (100) among all the schools in the country.

Using the same sample, conduct the appropriate hypothesis test with $\alpha = 0.05$.

Step 1: Assumptions: the type of data is quantitative data; the sample is generated by random sampling; the population is distributed normally.

Step 2: state hypothesis: Null hypothesis: the average student IQ in the school is lower than or equal to the average IQ score (100), ($H_0 \leq 100$)

Alternative hypothesis: $H_a > 100$

```

1 SAMMEAN ## sample mean
2 IQMEAN <- 100 ## average IQ score among all the schools in the country
3 SAMSD ## standard deviation of the sample
4 SAMSE <- SAMSD/sqrt(n)    ##sample standard error
5
6 H0 <- IQMEAN

```

Step 3: Calculate a test statistics:

```

1 TS <- (SAMMEAN - H0)/SAMSE
2 TS

```

report TS = -0.1191488

Step 4: P-value:

```

1 ZS <- (SAMMEAN-H0)/2.618    ## Z-score
2 ZS
3 PV <- 1-pnorm(-abs(-0.596))
4 PV

```

Report P-value = 0.72, larger than 0.05

or

```

1 t.test(y, mu = 100, alternative = "greater")

```

Report P-value = 0.72, larger than 0.05

Step 5: Draw a conclusion

The p-value is higher than 0.05, so we cannot deny the null hypothesis that the average IQ of students in this school is lower than or equal to the average IQ of students in all schools in this country.

Question 2 (50 points): Political Economy

Researchers are curious about what affects the amount of money communities spend on addressing homelessness. The following variables constitute our data set about social welfare expenditures in the USA.

State	50 states in US
Y	per capita expenditure on shelters/housing assistance in state
X1	per capita personal income in state
X2	Number of residents per 100,000 that are "financially insecure" in state
X3	Number of people per thousand residing in urban areas in state
Region	1=Northeast, 2= North Central, 3= South, 4=West

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

```
1 #####
2 # Problem 2
3 #####
4
5 expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI_Fall2021/main/datasets/expenditure.txt", header=T)
6 ## this link does not work. I also tried read.csv() but all columns
7 # fall into one, so I loaded the local file
8
9
10 expenditure <- read.table("C:/Users/Caesar/Documents/GitHub/StatsI_Fall2022/datasets/expenditure.txt")
11
12
13 # Researchers are curious about what affects the amount of money communities
14 # spend on addressing homelessness. The following variables constitute our
15 # data set about social welfare expenditures in the USA.
16
17
18 ## Check the structure, class and type of data
19 str(expenditure) ## The Structure of data is a data frame of character vectors
20 class(expenditure) ## The Class of data is data frame
21 typeof(expenditure) ## The type of data is list.
22 head(expenditure) ## The first row should be removed
23
24 ## change first row into headers
25 names(expenditure) <- expenditure[1, ]
26 expenditure <- expenditure[-1, ]
27 head(expenditure)
28
29
30 ## save as csv. file
31 write.csv(expenditure, "C:/Users/Caesar/Documents/GitHub/StatsI_Fall2022/
```

```

32     problemSets/PS01/My-Answers/exenditure.csv")
33 expenditure_csv <- read.csv("C:/Users/Caesar/Documents/GitHub/StatsI_Fall2022/
    problemSets/PS01/My-Answers/exenditure.csv")

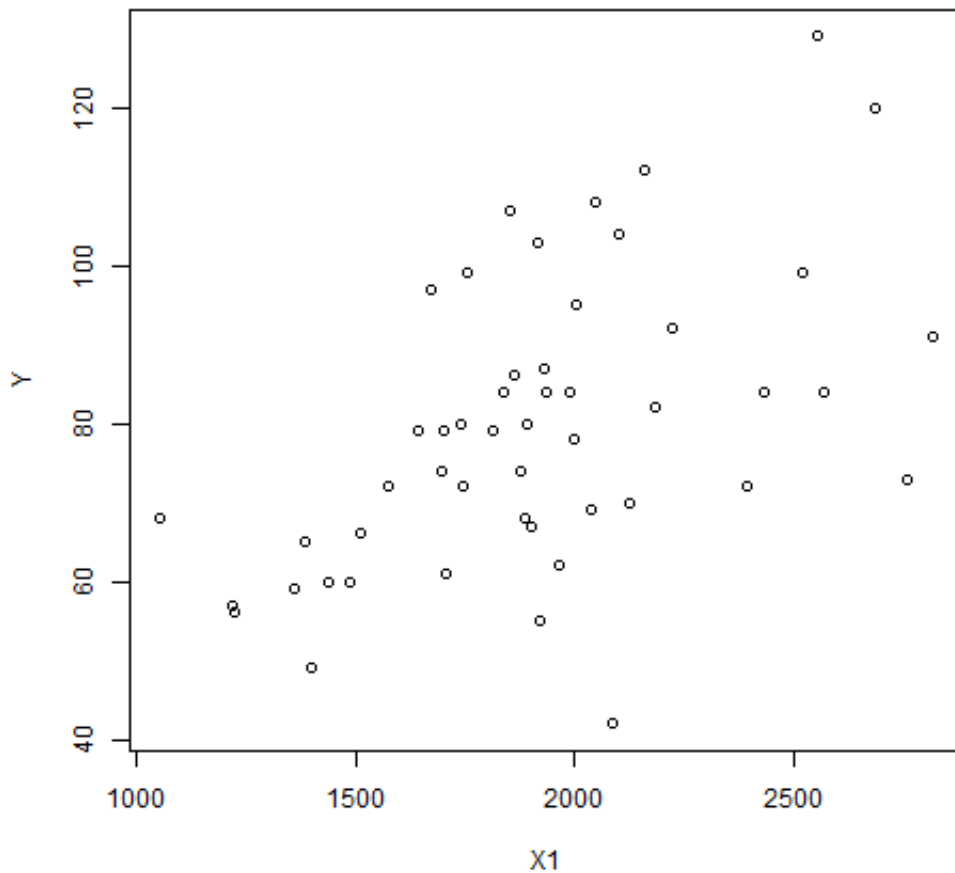
```

- 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)?

```

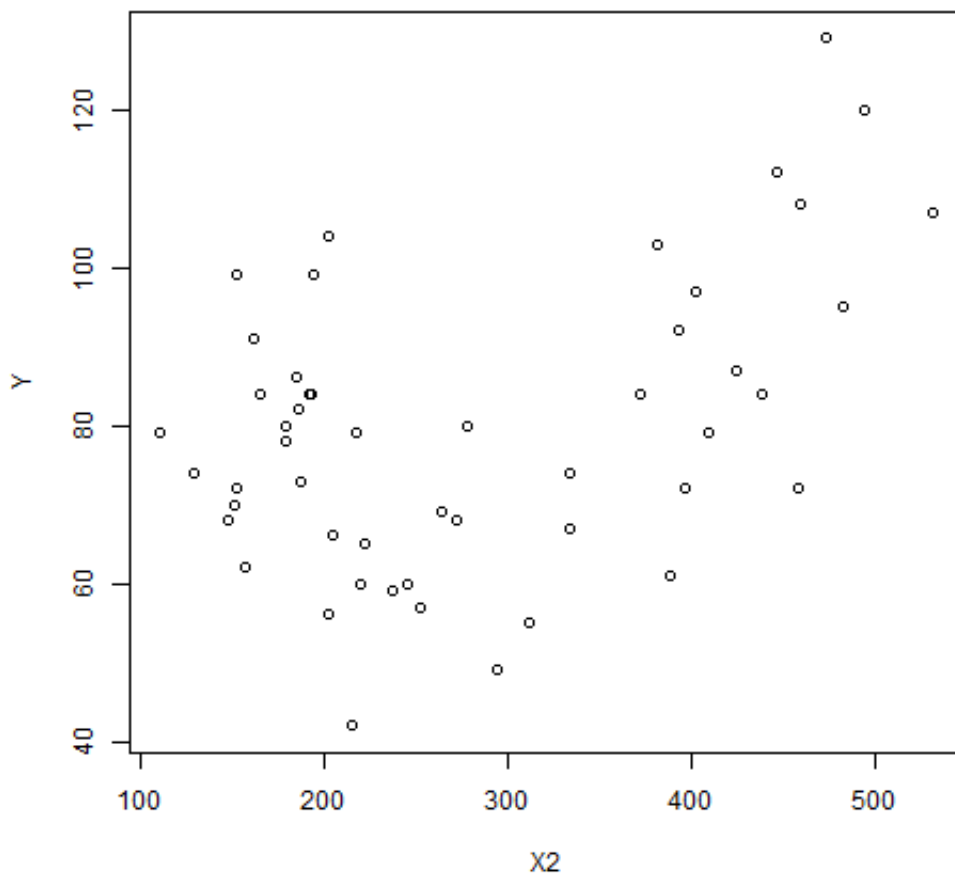
1 png("Y ~ X1.png")
2 plot(Y ~ X1, data = expenditure_csv)
3 # Preliminary analysis of the graph: When X1 increases, we expect to see
  an
4 # increase in Y at the same time.
5 dev.off()
6 REGYX1 <- lm(Y ~ X1, data = expenditure_csv)
7 summary(REGYX1)

```



Preliminary analysis of the graph: When X1 increases, we expect to see an increase in Y at the same time. The coefficient of X1 is 0.025, the intercept is 32.546, which means when X1, per capita personal income in state increases by 1 US dollar, we expect to see an increase in Y, per capita expenditure on shelters/housing assistance in state, by 0.025 US dollar; When per capita personal income in state is 0 dollar, the per capita expenditure on shelters/housing assistance in state is 32.546 US dollar. The p value of the coefficient of X1 is 7.08e-05, which is smaller than 0.001. Therefore, We can reject the null hypothesis that there is no correlation between Y and X1 at 0.001 level.

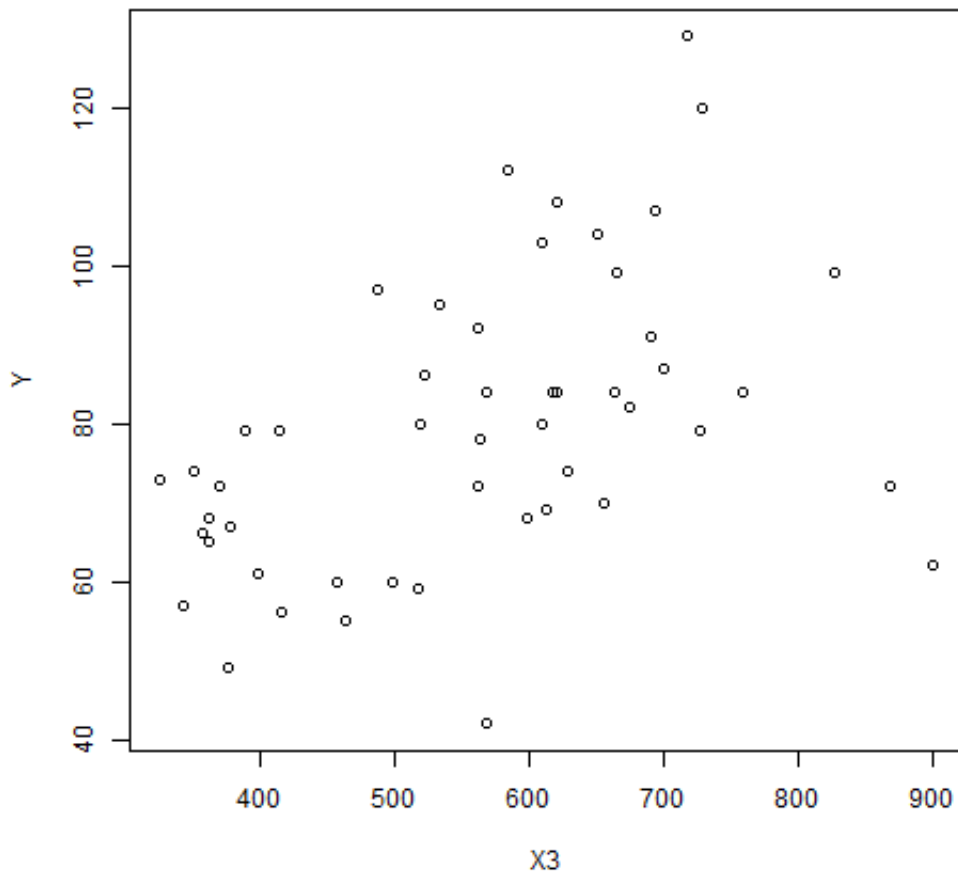
```
1 png("Y ~ X2.png")
2 plot(Y ~ X2, data = expenditure_csv)
3 # Preliminary analysis of the graph: When X2 increases , we expect to see
  an
4 # increase in Y at the same time.
5 dev.off()
6 REGYX2 <- lm(Y ~ X2, data = expenditure_csv)
7 summary(REGYX2)
```



Preliminary analysis of the graph: When X2 increases, we expect to see an increase in Y at the same time.

The coefficient of X2 is 0.070, the intercept is 59.761, which means when X2, the number of residents per 100,000 that are "financially insecure" in state increases by 1, we expect to see an increase in Y, per capita expenditure on shelters/housing assistance in state, by 0.070 US dollar; When Number of residents per 100,000 that are "financially insecure" in state is 0, the per capita expenditure on shelters/housing assistance in state is 59.761 US dollar. The p value of the coefficient of X2 is 0.001, which is smaller than 0.01. Therefore, We can reject the null hypothesis that there is no correlation between Y and X2 at 0.01 level.

```
1 png("Y ~ X3.png")
2 plot(Y ~ X3, data = expenditure_csv)
3 # Preliminary analysis of the graph: When X3 increases , we expect to see
  Y
4 # increase at the same time.
5 dev.off()
6 REGYX3 <- lm(Y ~ X3, data = expenditure_csv)
7 summary(REGYX3)
```



Preliminary analysis of the graph: When X3 increases, we expect to see Y increase at the same time.

The coefficient of X3 is 0.059, the intercept is 46,306, which means when X3, the number of people per thousand residing in urban areas in state increases by 1, we expect to see an increase in Y, per capita expenditure on shelters/housing assistance in state, by 0.059 US dollar; When Number of people per thousand residing in urban areas in state is 0, the per capita expenditure on shelters/housing assistance in state is 46.306 US dollar. The p value of the coefficient of X3 is 0.000695, which is smaller than 0.001. Therefore, We can reject the null hypothesis that there is no correlation between Y and X3 at 0.001 level.

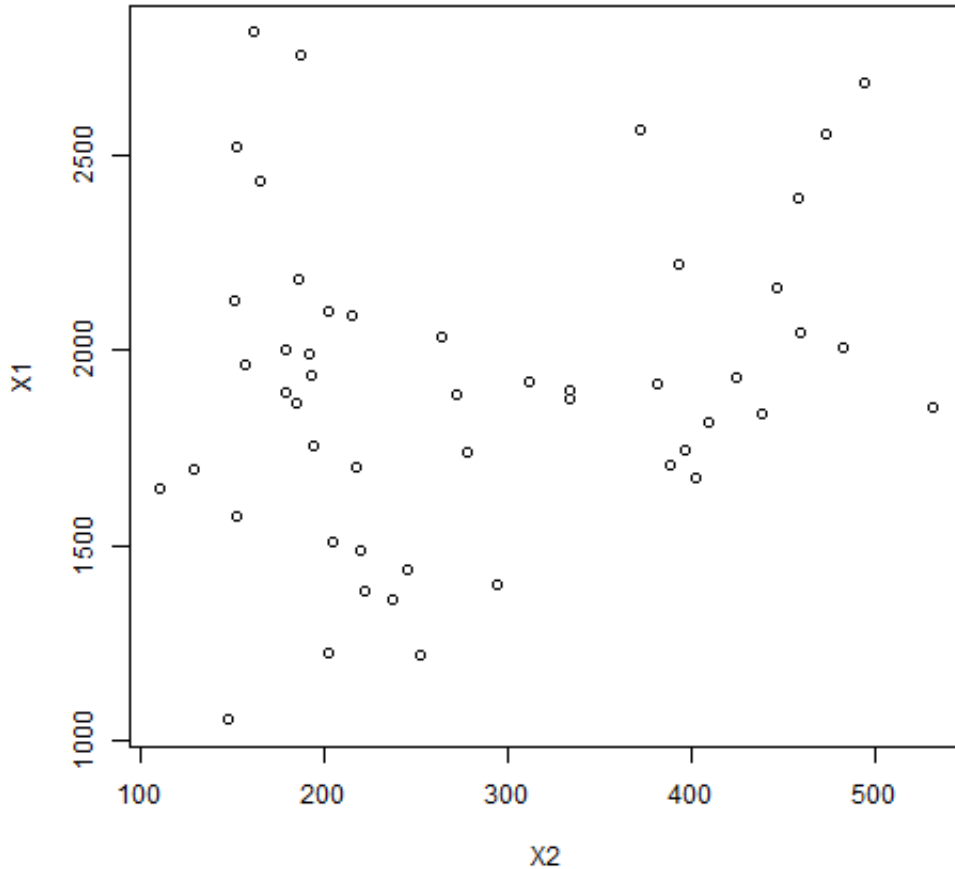
```

1 png("X1 ~ X2.png")
2 plot(X1 ~ X2, data = expenditure_csv)
3 # Preliminary analysis of the graph: there is no obvious correlation
4 # between X1 and X2.
5 dev.off()
6 REGX1X2 <- lm(X1 ~ X2, data = expenditure_csv)

```



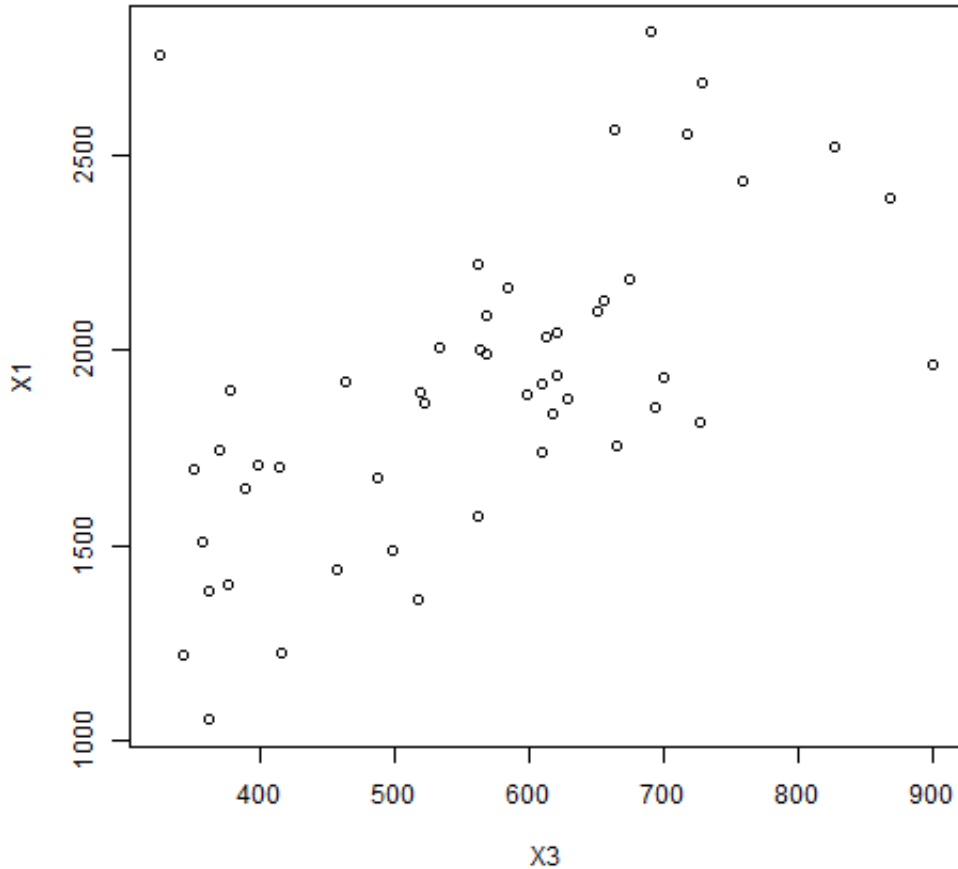
```
7 summary(REGX1X2)
```



Preliminary analysis of the graph: there is no obvious correlation between X1 and X2. The coefficient of X2 is 0.696, the intercept is 1715.655, which means when X2, number of residents per 100,000 that are "financially insecure" in state increases by 1, we expect to see an increase in X1, per capita personal income in state, by 0.696 US dollar; When the number of residents per 100,000 that are "financially insecure" in state is 0, the per capita personal income in state is 1715.655 US dollar. The p value of the coefficient of X2 is 0.152, which is larger than 0.05. Therefore, We cannot reject the null hypothesis that there is no correlation between X1 and X2.

```
1 png("X1 ~ X3.png")
2 plot(X1 ~ X3, data = expenditure_csv)
3 # Preliminary analysis of the graph: when X3 increases, X1 is expected to
4 # increase at the same time.
5 dev.off()
6 REGX1X3 <- lm(X1 ~ X3, data = expenditure_csv)
```

```
7 summary(REGX1X3)
```



Preliminary analysis of the graph: when X3 increases, X1 is expected to increase at the same time.

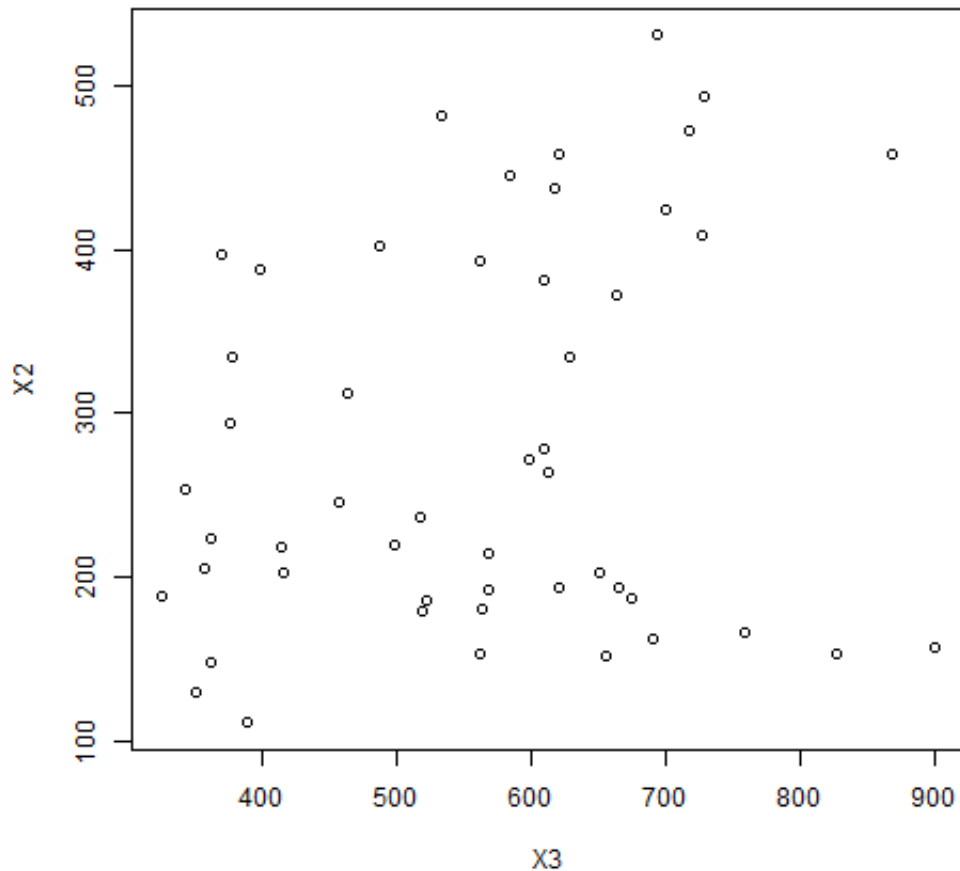
The coefficient of X3 is 1.643, the intercept is 988.947, which means when X3, the number of people per thousand residing in urban areas in state increases by 1, we expect to see an increase in X1, per capita personal income in state, by 1.643 US dollar; When the Number of people per thousand residing in urban areas in state is 0, the per capita personal income in state is 988.947 US dollar. The p value of the coefficient of X3 is 5.13e-06, which is smaller than 0.001. Therefore, We can reject the null hypothesis that there is no correlation between X1 and X2 at 0.001 level.

```
1 png("X2 ~ X3.png")
2 plot(X2 ~ X3, data = expenditure_csv)
3 # Preliminary analysis of the graph: There is no obvious correlation
4 # between X2 and X3.
5 dev.off()
```

```

6 REGX2X3 <- lm(X2 ~ X3, data = expenditure_csv)
7 summary(REGX2X3)

```



Preliminary analysis of the graph: There is no obvious correlation between X2 and X3.

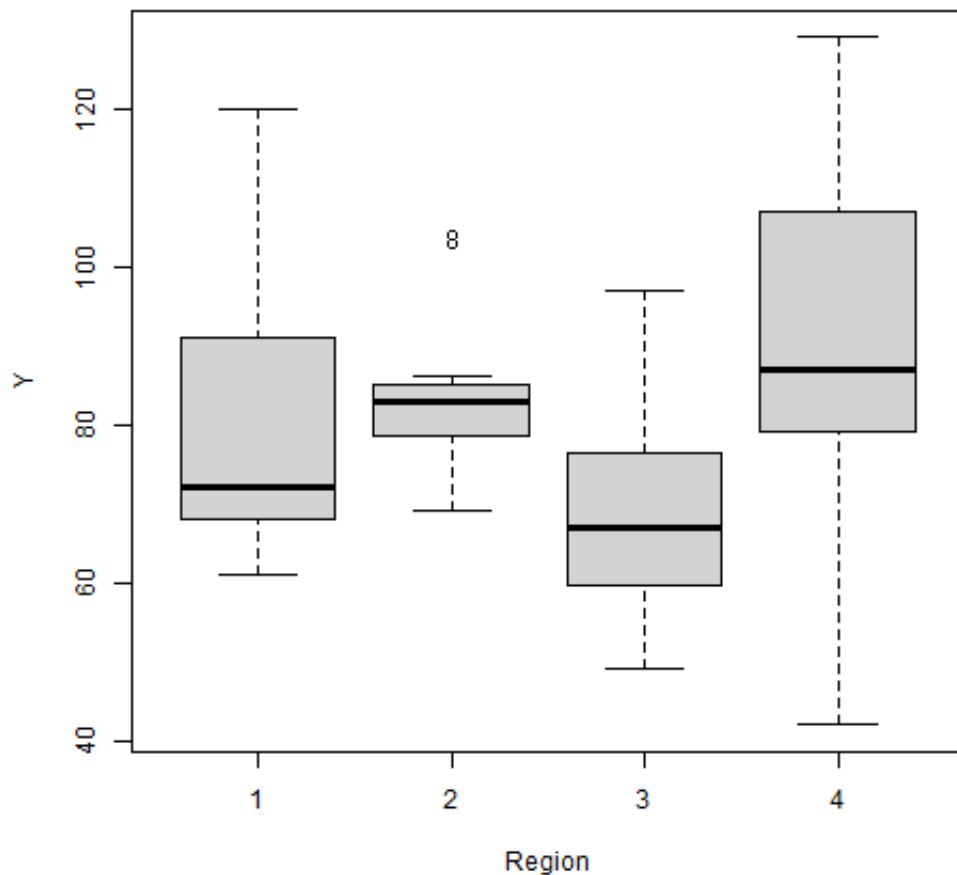
The coefficient of X3 is 0.180, the intercept is 180.609, which means when X3, the number of people per thousand residing in urban areas in state increases by 1, we expect to see an increase in X2, Number of residents per 100,000 that are "financially insecure" in state by 0.180; When the number of people per thousand residing in urban areas in state is 0, the number of residents per 100,000 that are "financially insecure" in state is 180.609. The p value of the coefficient of X3 is 0.123, which is larger than 0.05. Therefore, We cannot reject the null hypothesis that there is no correlation between X2 and X3.

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

```

1 png("Y ~ Region.png")
2 boxplot(Y ~ Region, data = expenditure_csv)
3 ## on average, West has the highest per capita expenditure on housing
  assistance.
4 dev.off()

```



On average, West has the highest per capita expenditure on housing assistance.

- 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.

```

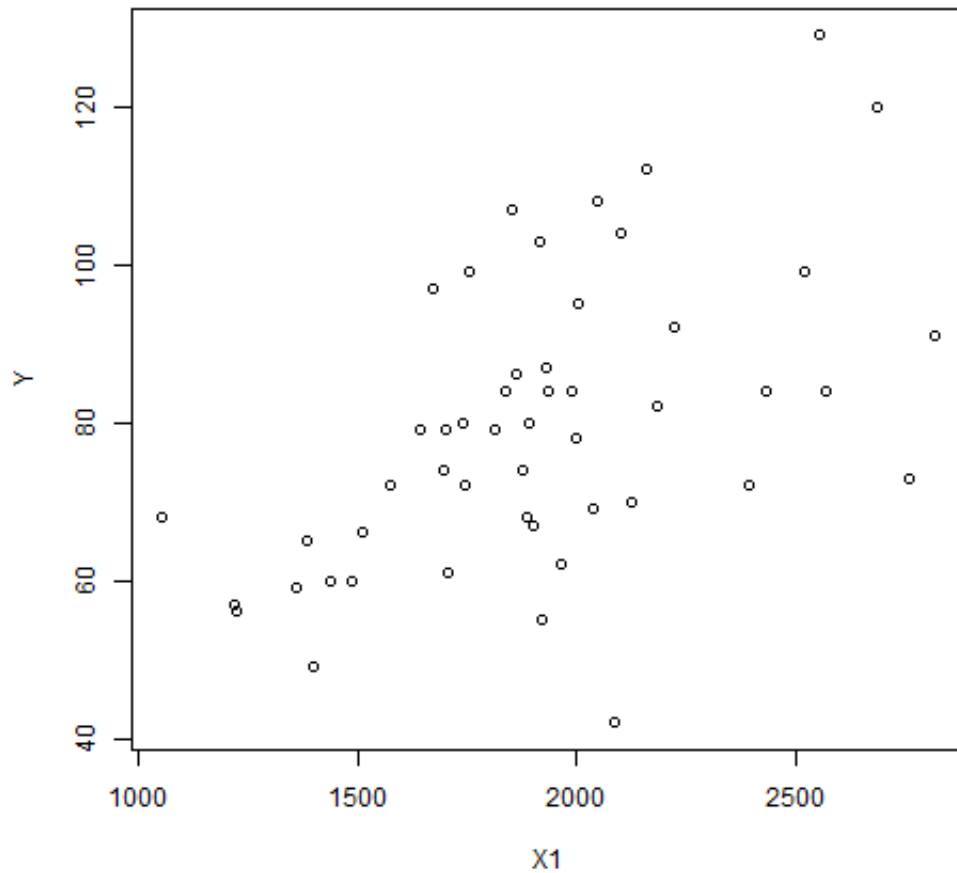
1 png("Y ~ X1.png")
2 plot(Y ~ X1, data = expenditure_csv)
3 # Preliminary analysis of the graph: When X1 increases, we expect to see
  an
4 # increase in Y at the same time.

```

```

5 dev.off()
6 REGYX1 <- lm(Y ~ X1, data = expenditure_csv)
7 summary(REGYX1)

```



Preliminary analysis of the graph: When X1 increases, we expect to see an increase in Y at the same time.

The coefficient of X1 is 0.025, the intercept is 32.546, which means when X1, per capita personal income in state increases by 1 US dollar, we expect to see an increase in Y, per capita expenditure on shelters/housing assistance in state, by 0.025 US dollar; When per capita personal income in state is 0 dollar, the per capita expenditure on shelters/housing assistance in state is 32.546 US dollar. The p value of the coefficient of X1 is 7.08e-05, which is smaller than 0.001. Therefore, We can reject the null hypothesis that there is no correlation between Y and X1 at 0.001 level.

```

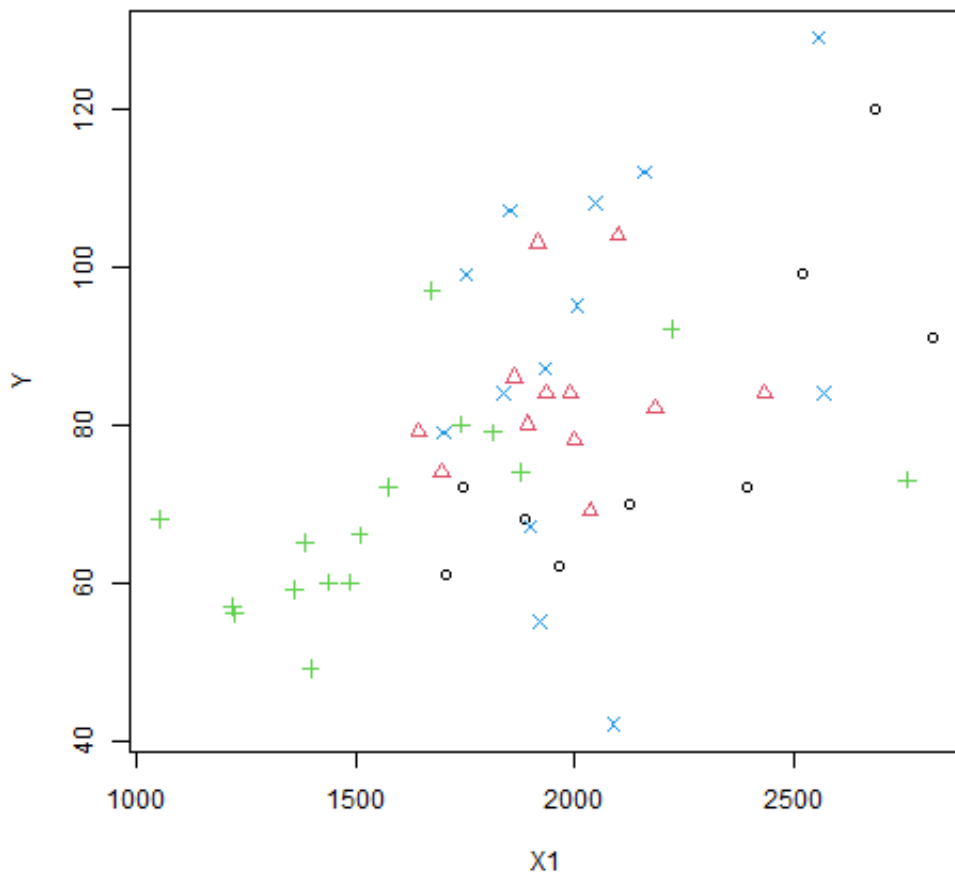
1 png("Y ~ X1, colour, symbol.png")
2 plot(Y ~ X1, data = expenditure_csv, col = Region, pch = Region)
3

```

```

4 dev.off()
5 REGYX1 <- lm(Y ~ X1, data = expenditure_csv)
6 summary(REGYX1)
7 REGYX1REG <- lm(Y ~ X1 + Region, data = expenditure_csv)
8 summary(REGYX1REG)

```



The coefficient of X1 is 0.027, the coefficient of Region is 3.333, which means when X1, per capita personal income in state increases by 1 US dollar, we expect to see an increase in Y, per capita expenditure on shelters/housing assistance in state, by 0.027 US dollar; when the region is different, we expect to see an average difference in the per capita expenditure on shelters/housing assistance in state of 3.333 US dollar. The p value of the coefficient of X1 is 2.77e-05, which is smaller than 0.001. Therefore, We can reject the null hypothesis that there is no correlation between Y and X1 at 0.001 level.

The p value of the coefficient of Region is 0.128, which is larger than 0.05. Therefore, we cannot reject the null hypothesis that there is no correlation between Y and Region.