

Problem Set 1

Applied Stats/Quant Methods 1

Due: September 30, 2024

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.
- This problem set is due before 23:59 on Monday September 30, 2024. No late assignments will be accepted.

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

```
1
2 y <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113,
3       112, 98, 80, 97, 95, 111, 114, 89, 95, 126, 98)
4 # Section1: Calculate the 90% Confidence Interval
5
6 # Given the sample size is less than 30(n=25),so we can use t-test to
7   calculate the confidence interval
8 # Sample size
9 n <- length(y)
10
11 # Set the significance level and get the t-critical value
12 alpha <- 0.10
13 t_critical <- qt(0.05, df=n-1)
```

```

13
14 # Calculate the standard deviation of the sample
15 sample_sd <- sd(y)
16
17 # This value reflects the uncertainty around the sample mean
18 margin_of_error <- t_critical * (sd(y) / sqrt(n))
19
20 #Calculate the confidence interval limits
21 CI_lower <- mean(y) - margin_of_error
22 CI_upper <- mean(y) + margin_of_error
23
24 #print confidence interval
25 print(c(CI_lower, mean(y), CI_upper))
26
27 # The calculated 90% confidence interval is approximately [93.95, 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 $\alpha = 0.05$.

```

1 # Section2:Conduct the Hypothesis Test
2
3 #Perform a one-sample t-test to test whether the hypothesized mean is
  significantly greater than 100
4 # The null hypothesis (H0) is that the true mean IQ is equal to 100.
5 # The alternative hypothesis (H1) is that the true mean IQ is greater
  than 100.
6
7 #Perform the t-test
8 t_test_result <- t.test(y, mu = 100, alternative = "greater", conf.level
  = 0.95)
9
10 #Print the result of the t-test
11 print(t_test_result)
12
13 # The output includes:
14 # t-value: -0.59574
15 # p-value: 0.7215, which is quite large and suggests we fail to reject
  the null hypothesis.
16 # Alternative hypothesis: true mean is greater than 100.
17 # 95% Confidence Interval: [93.96, Inf] – The lower bound is 93.96, with
  no upper bound (as it's a one-tailed test).
18 # Sample mean: 98.44, which is less than 100
19 # aligning with the test's result that the average IQ is not
  significantly greater than 100.

```

The average IQ of students in her school is higher than the average IQ

Part I

title

Question 2: 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.

- 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)?

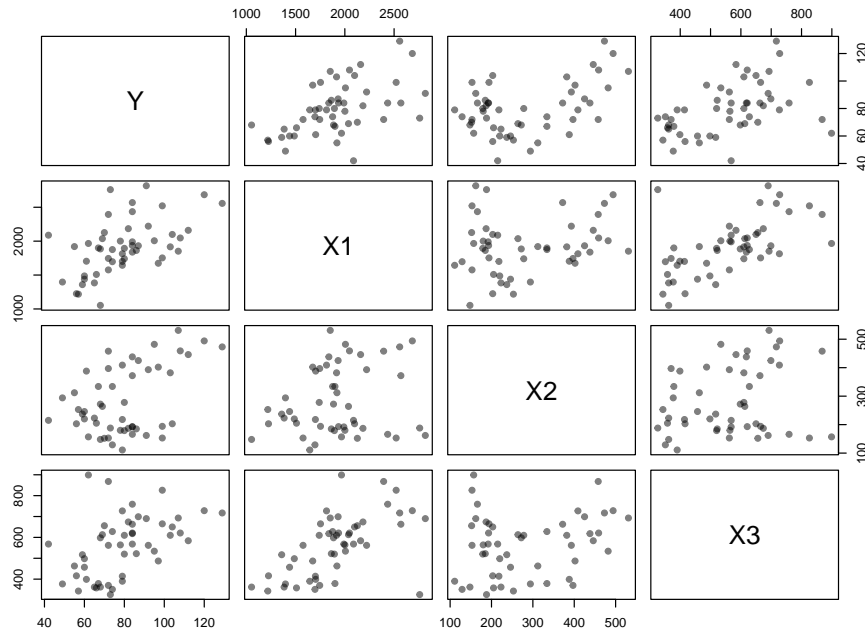
```
1 # Load required libraries
2 library(ggplot2)
3 library(dplyr)
4 library(corrplot)
5
6 # Import the data
7 expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/
8   StatsI_Fall2024/main/datasets/expenditure.txt", header = TRUE)
9 head(expenditure)
10
11 # Rename columns for clarity
12 colnames(expenditure) <- c("State", "Y", "X1", "X2", "X3", "Region")
13
14 # Convert Region to factor
15 expenditure$Region <- factor(expenditure$Region, levels = 1:4, labels = c("
16   Northeast", "North Central", "South", "West"))
17
18 # Create correlation matrix
19 cor_matrix <- cor(expenditure[, c("Y", "X1", "X2", "X3")])
20 print(cor_matrix)
21
22 # section 1: Plot relationships among Y, X1, X2, and X3
23 # Create scatter plots for all bivariate relationships
24 pdf(file="scatter_Matrix.pdf", width=8, height=6)
```

```

24 # Plot the pairwise scatterplots between Y, X1, X2, X3
25 pairs(expenditure[, c("Y", "X1", "X2", "X3")], pch = 19, col = alpha("
    black", 0.5))
26 dev.off()
27
28 # The correlation among Y, X1, X2, and X3
29 # Y and X1 show a moderate positive correlation (~0.53), suggesting that
    states with higher personal income tend to spend more on housing
    assistance.
30 # Y and X2 have a moderate positive correlation (~0.45), indicating a
    link between the number of financially insecure residents and housing
    expenditure.
31 # Y and X3 also show a moderate positive correlation (~0.46), implying
    that urbanization (X3) has some relationship with higher housing
    assistance spending.
32 # The relationships between X1 and X2, X1 and X3, and X2 and X3 show
    varying degrees of correlation.

```

Figure 1: scatter Matrix



The scatter plots show a moderate positive relationship between per capita expenditure on housing assistance (Y) and each of the predictors: personal income (X1), financially insecure residents (X2), and urban population percentage (X3). This suggests that states with higher incomes, larger urban populations, and more financially vulnerable residents tend to spend more on addressing homelessness.

Figure 2: Y X1

Per Capita Expenditure vs Per Capita Personal Income

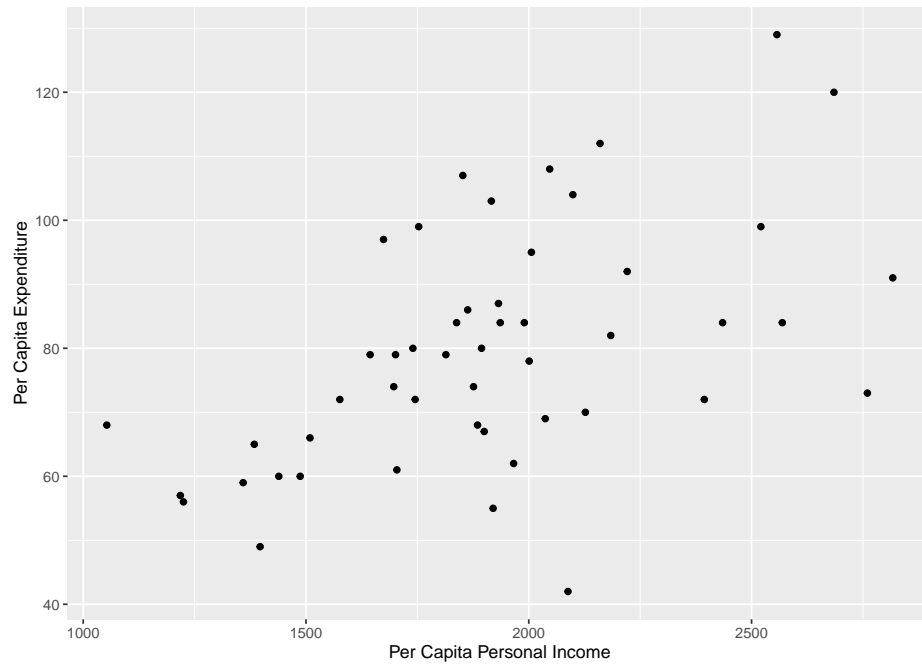


Figure 3: X2 Y

Per Capita Expenditure vs Financially Insecure Residents

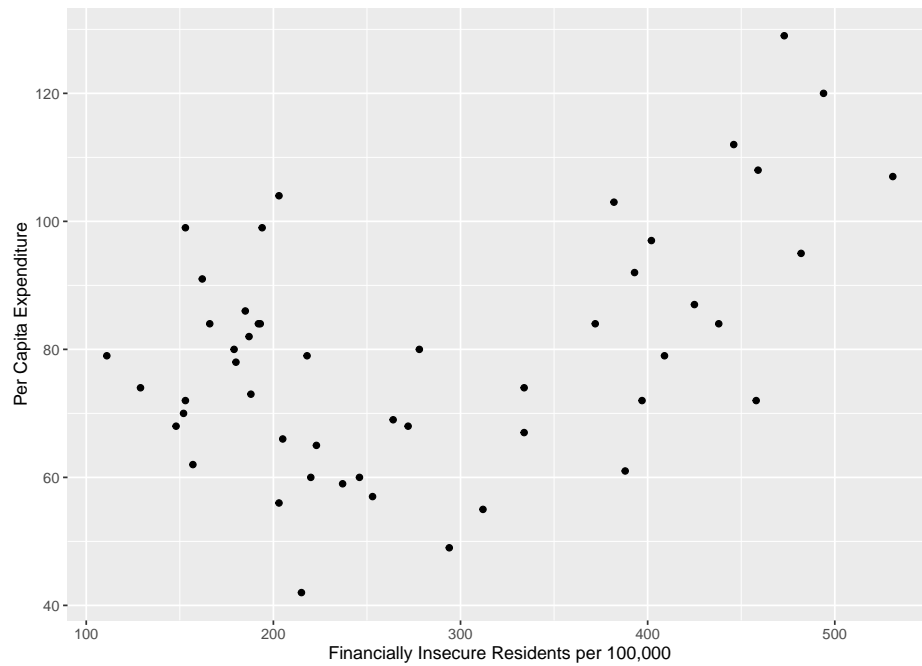


Figure 4: X3 Y

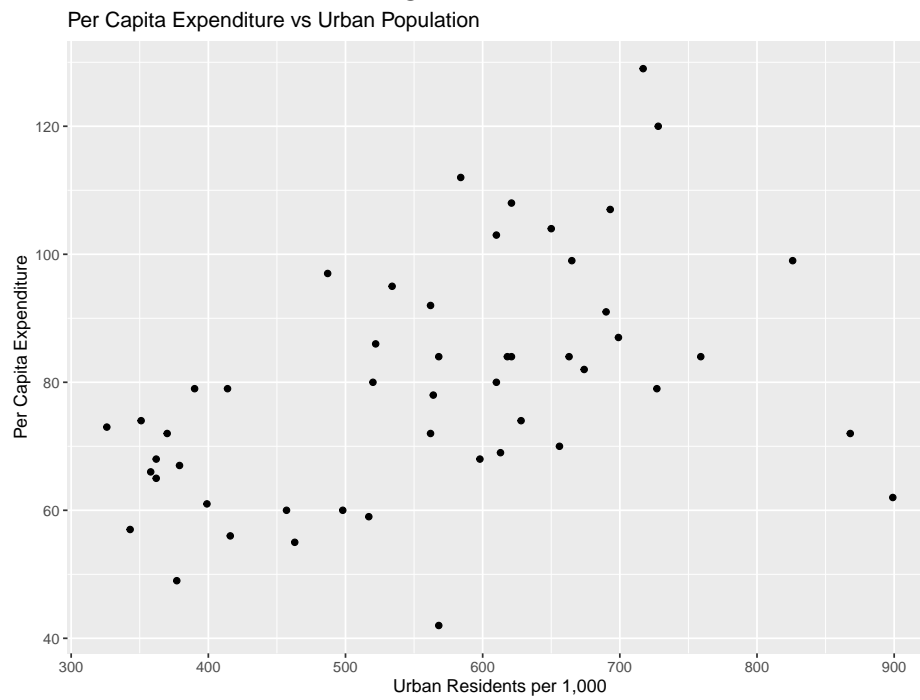


Figure 5: X1 X2

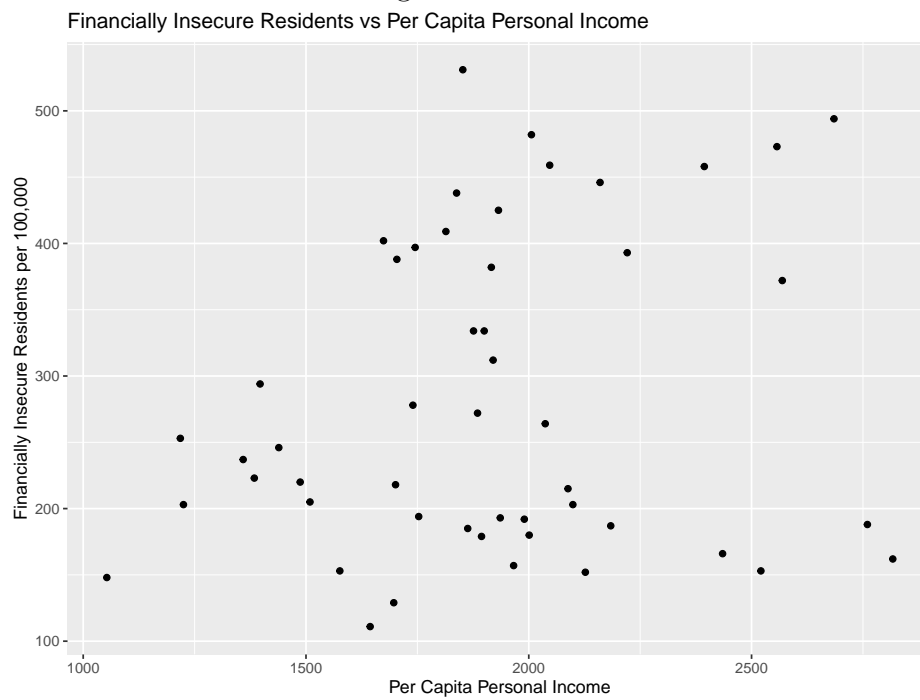
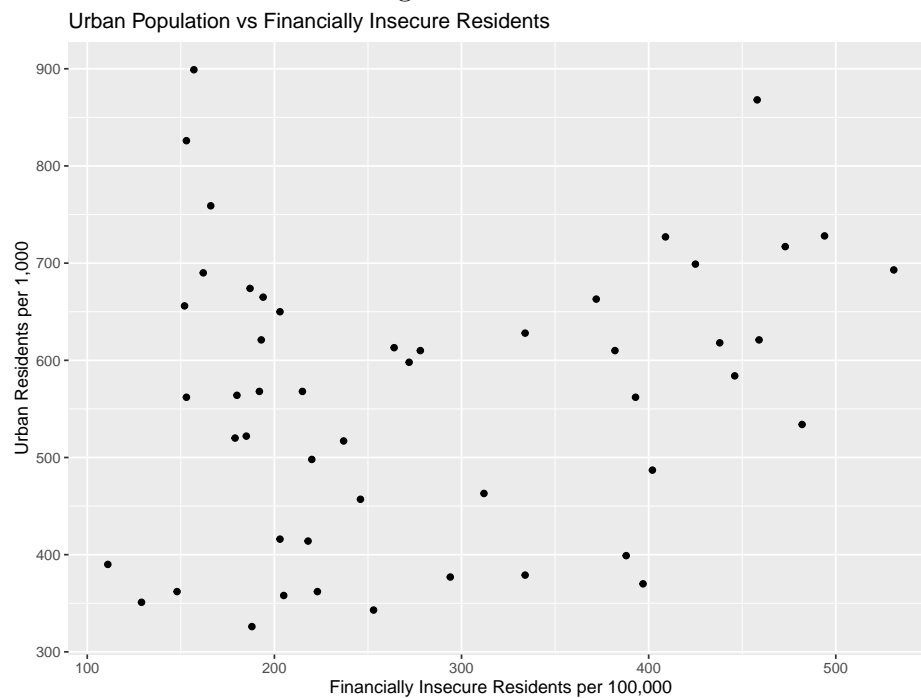


Figure 6: X1 X3

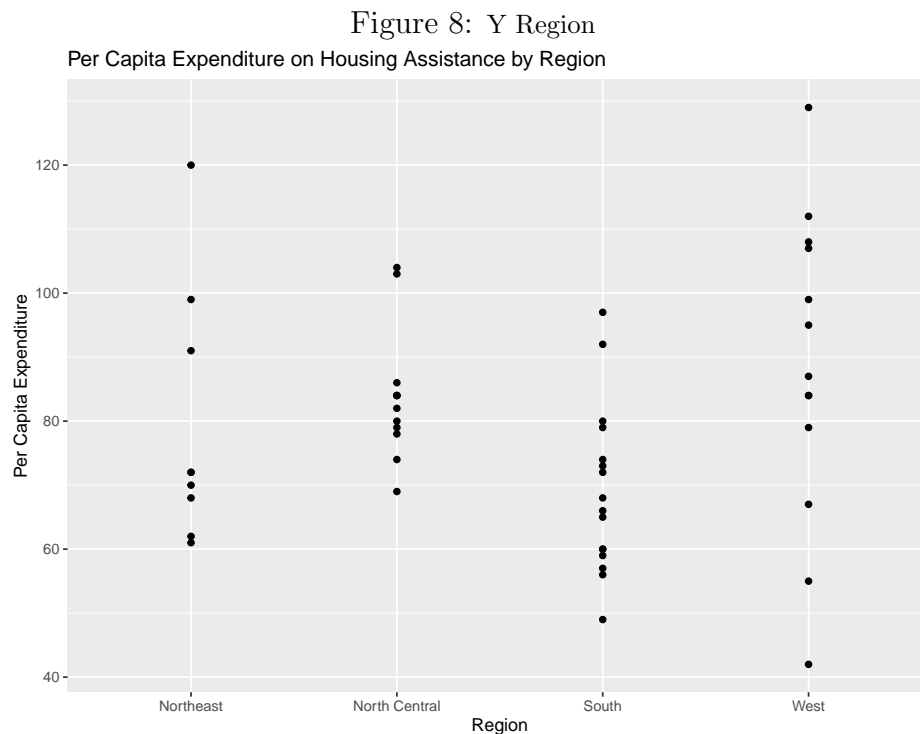


Figure 7: X2 X3



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

```
1 #section 2: Plot the relationship between Y and Region
2 Y_vs_Region <- ggplot(expenditure, aes(x = Region, y = Y)) +
3   geom_point() +
4   labs(title = "Per Capita Expenditure on Housing Assistance by Region",
5         x = "Region", y = "Per Capita Expenditure")
6 ggsave("Y_VS_Region.pdf", plot=Y_vs_Region, width=8, height=6, units="in")
7
8 # Based on the plot, Region 1 (Northeast) has the highest average 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.

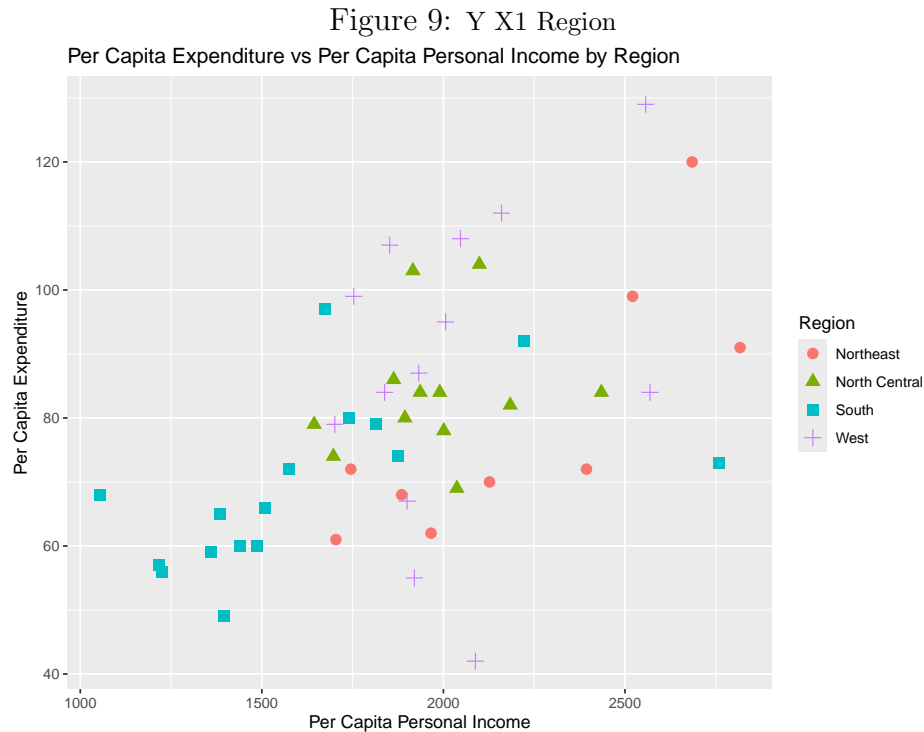
```
1 # Section3: Plot the relationship between Y and X1
2 Y_vs_X1 <- ggplot(expenditure, aes(x = X1, y = Y)) +
3   geom_point() +
4   labs(title = "Per Capita Expenditure vs Per Capita Personal Income",
5         x = "Per Capita Personal Income", y = "Per Capita Expenditure")
6 ggsave("Y_vs_X1.pdf", plot=Y_vs_X1, width=8, height=6, units="in")
7
8 # The scatterplot shows a positive linear relationship between Y and X1 (
   per capita income).
```



```

9 # Higher personal income in a state tends to be associated with higher
  spending on housing assistance.
10
11 # Reproduce the Y vs X1 scatterplot, adding Region as a distinguishing
  factor
12 Y_X1_Region <- ggplot(expenditure, aes(x = X1, y = Y, color = Region,
  shape = Region)) +
13   geom_point(size = 3) +
14   labs(title = "Per Capita Expenditure vs Per Capita Personal Income by
  Region",
15         x = "Per Capita Personal Income", y = "Per Capita Expenditure")
16 ggsave("Y_X1_Region.pdf", plot=Y_X1_Region, width=8,height=6,units="in")
17
18 # The plot shows that different regions have different expenditure trends
  with respect to personal income.
19 # Region 1 (Northeast) tends to have higher spending across the income
  spectrum, while Region 3 (South) generally shows lower expenditure.

```



The scatter plot of Y (per capita expenditure on housing assistance), X1 (personal income), and Region shows a positive relationship between Y and X1, with higher-income states spending more on housing assistance. However, regional differences are clear, with the Northeast and West regions generally spending more at similar income levels compared to the South and North Central. This suggests that both personal income and regional factors play significant roles in determining housing assistance expenditures.