

Problem 1

Yajie Dong 23339234

September 29, 2023

Question 1: Education

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

Step 1: Calculate the sample mean

```
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 y_bar <- mean(y)
3 y_bar
```

$$\bar{y} = 98.44$$

Step 2: Calculate the Sample Standard Deviation (S)

```
1 S <- sd(y)
2 S
```

$$S = \sqrt{\frac{1}{24} \sum_{i=1}^{25} (y_i - 98.44)^2} \approx 13.09$$

Step 3: Find the Z-Score

```
1 confidence_level <- 0.9
2 alpha <- 1 - confidence_level
3 z_score <- qnorm(1 - alpha / 2)
4 z_score
```

Z-score is approximately 1.645

Step 4: Calculate the Margin of Error (ME)

```
1 ME <- z_score * (S / sqrt(n))  
2 ME
```

$$ME = z \times \frac{S}{\sqrt{n}} \approx 1.645 \times \frac{13.09}{5} \approx 4.31$$

Step 5: Calculate the confidence Interval

```
1 CI_lower <- y_bar - ME  
2 CI_upper <- y_bar + ME  
3 cat("90% Confidence Interval for Average Student IQ: [", CI_lower, "  
  ",",", CI_upper, "]\n")
```

Lower limit: $\bar{y} - ME \approx 98.44 - 4.31 \approx 94.13$

Upper limit: $\bar{y} + ME \approx 98.44 + 4.31 \approx 102.75$

Therefore, the 90% confidence interval for the average student IQ in the school is approximately [94.13,102.75]

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.

Step 1: State Hypotheses

-Null Hypothesis $H_0 : \mu = 100$ (the average IQ in the school is not different from 100)

-Alternative Hypothesis $H_a : \mu > 100$ (the average IQ in the school is greater than 100)

```
1 mu <- 100  
2 alpha <- 0.05
```

The one-sided test is used in this question because the school counselor is specifically interested in finding out whether the average IQ of students in her school is "higher" than the average IQ of 100 for the general population.

Step 2: Calculate Test Statistic

```
1 Z <- (y_bar - mu) / (S / sqrt(n))
2 Z
```

$$Z = \frac{\bar{y} - \mu}{\frac{S}{\sqrt{n}}} = \frac{98.44 - 100}{\frac{13.09}{\sqrt{25}}} \approx -0.60$$

Step 3: Find the critical value

```
1 Z_alpha <- qnorm(1 - alpha)
2 Z_alpha
```

$$Z_{\alpha} \approx 1.28$$

Step 4: Calculate P-value

```
1 p_value <- 1 - pnorm(Z)
2 p_value
```

$$P\text{-value} \approx 0.72$$

Step 5: Make a Decision.

```
1 if (Z > Z_alpha) {
2   decision <- "Reject Null Hypothesis"
3 } else {
4   decision <- "Fail to Reject Null Hypothesis"
5 }
6 decision
```

Neither the Z-test $Z < Z_{\alpha}$ or the P-value ($P\text{-value} > \alpha$) provide evidence to reject the null hypothesis.

Since we fail to reject the null hypothesis, we do not have enough evidence to suggest that the average IQ of the students in the school is greater than 100.

Question 2: Political Economy

1. Plot the relationships among Y, X1, X2, and X3? What are the correlations among them.

Step 1 : Import the data and summary

```
1 expenditure <- read.table("https://raw.githubusercontent.com/ASDS-
  TCD/StatsI_Fall2023/main/datasets/expenditure.txt", header =
  TRUE)
2
3 summary_stats <- summary(expenditure[, c("Y", "X1", "X2", "X3")])
4 cat("Dataset Summary:\n")
5 print(summary_stats)
```

Y	X1	X2	X3
Min. : 42.00	Min. :1053	Min. :111.0	Min. :326.0
1st Qu.: 67.25	1st Qu.:1698	1st Qu.:187.2	1st Qu.:426.2
Median : 79.00	Median :1897	Median :241.5	Median :568.0
Mean : 79.54	Mean :1912	Mean :281.8	Mean :561.7
3rd Qu.: 90.00	3rd Qu.:2096	3rd Qu.:391.8	3rd Qu.:661.2
Max. :129.00	Max. :2817	Max. :531.0	Max. :899.0

Step2: Calculate and print the correlation matrix

```
1 cor_matrix <- cor(expenditure[, c("Y", "X1", "X2", "X3")])
2 cat("Correlation Matrix:\n")
3 print(cor_matrix)
```

	Y	X1	X2	X3
Y	1.0000000	0.5317212	0.4482876	0.4636787
X1	0.5317212	1.0000000	0.2056101	0.5952504
X2	0.4482876	0.2056101	1.0000000	0.2210149
X3	0.4636787	0.5952504	0.2210149	1.0000000

Step 3 :Plot the relationships among Y, X1, X2, and X3

```
1
2 plot1 <- ggplot(expenditure, aes(x = X1, y = Y)) + geom_point() +
  geom_smooth(method = 'lm') + ggtitle("Y vs X1")
3 plot2 <- ggplot(expenditure, aes(x = X2, y = Y)) + geom_point() +
  geom_smooth(method = 'lm') + ggtitle("Y vs X2")
4 plot3 <- ggplot(expenditure, aes(x = X3, y = Y)) + geom_point() +
  geom_smooth(method = 'lm') + ggtitle("Y vs X3")
5 plot4 <- ggplot(expenditure, aes(x = X1, y = X2)) + geom_point() +
  geom_smooth(method = 'lm') + ggtitle("X1 vs X2")
6 plot5 <- ggplot(expenditure, aes(x = X1, y = X3)) + geom_point() +
  geom_smooth(method = 'lm') + ggtitle("X1 vs X3")
7 plot6 <- ggplot(expenditure, aes(x = X2, y = X3)) + geom_point() +
  geom_smooth(method = 'lm') + ggtitle("X2 vs X3")
8 print(plot1)
9 print(plot2)
10 print(plot3)
11 print(plot4)
12 print(plot5)
13 print(plot6)
```

Figure 1: Y and X1.

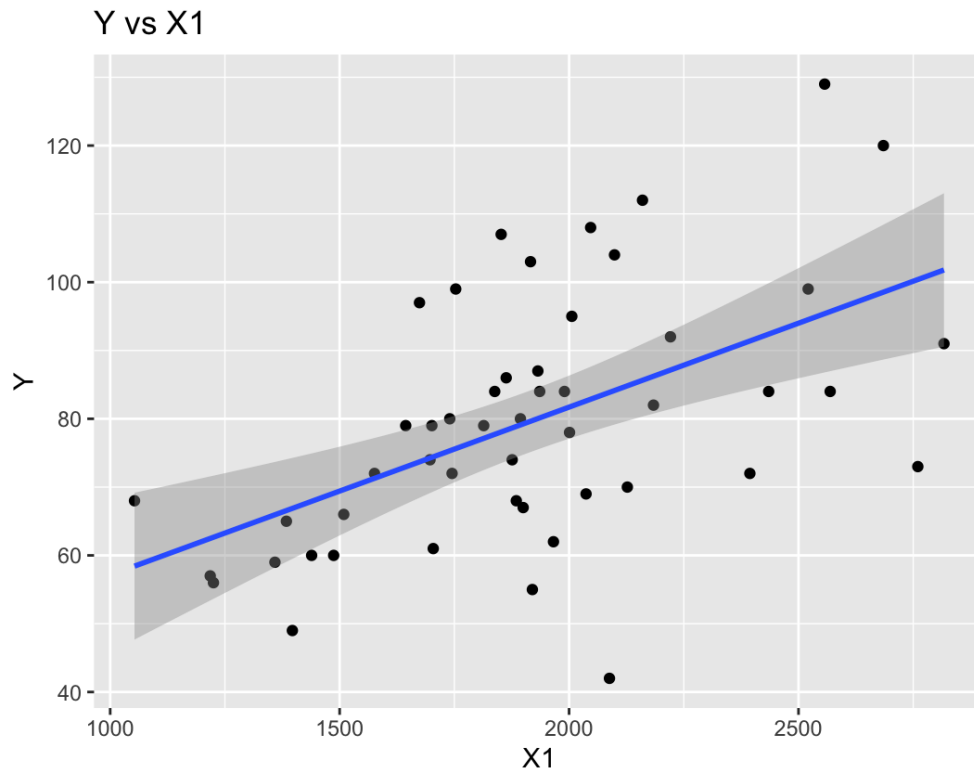


Figure 2: Y and X2.

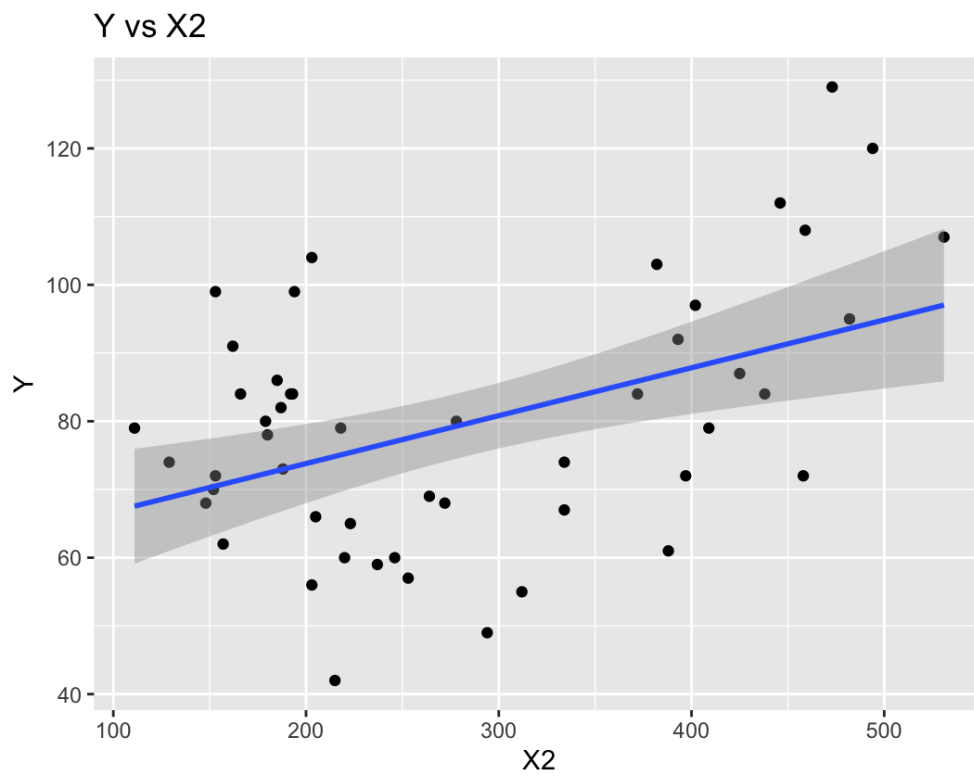


Figure 3: Y and X3.

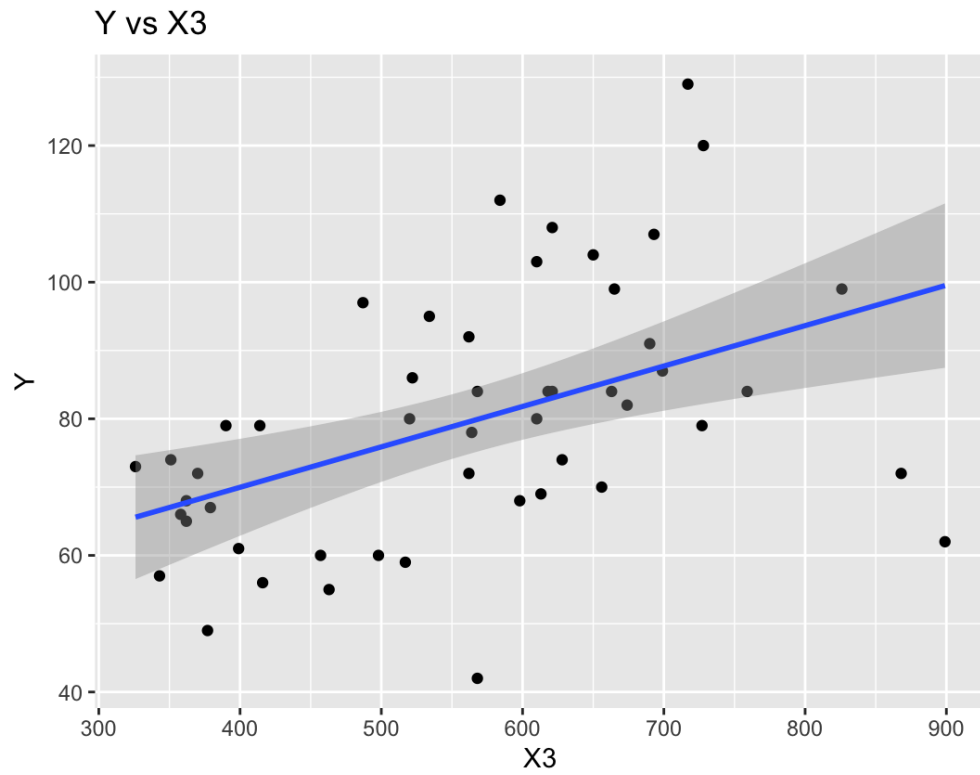


Figure 4: X1 and X2.

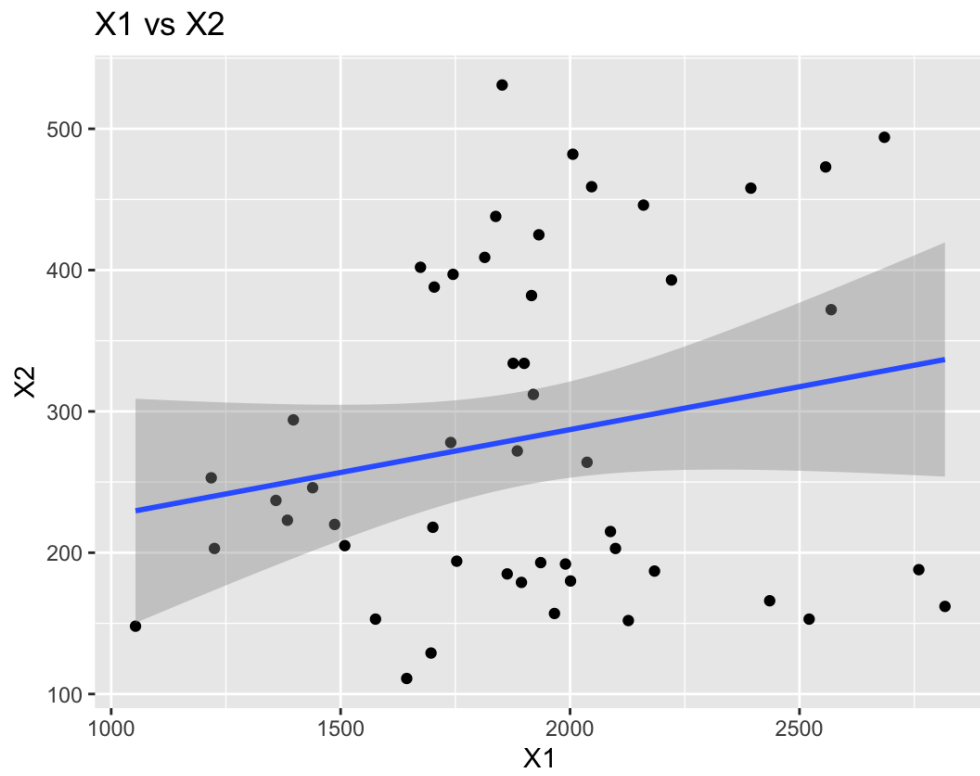


Figure 5: X_1 and X_3 .

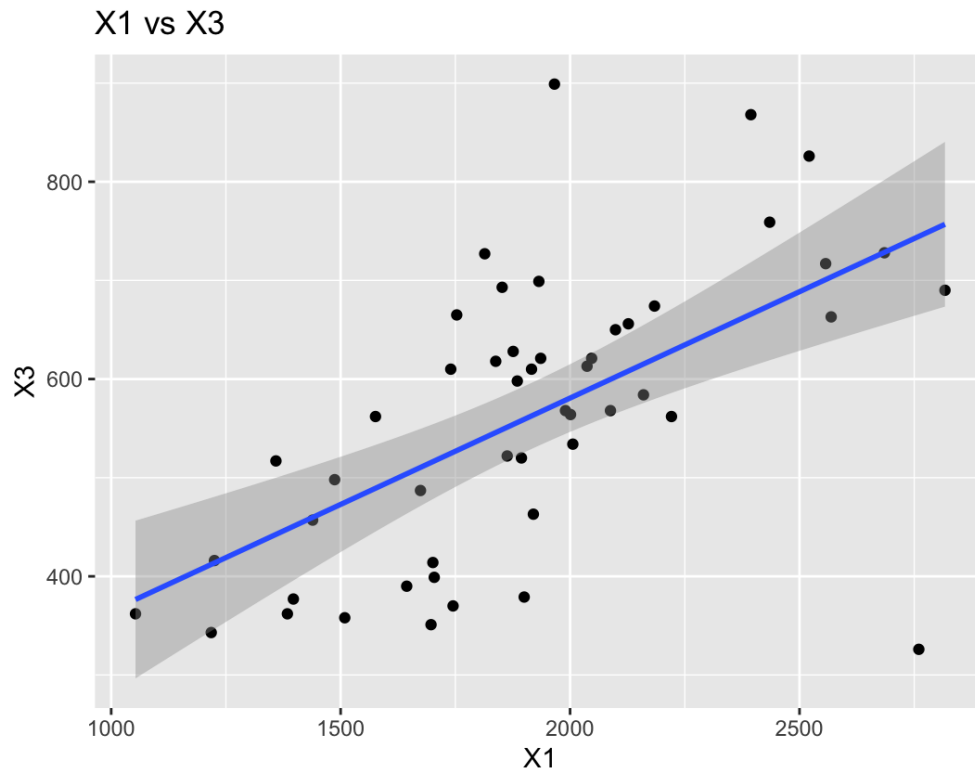
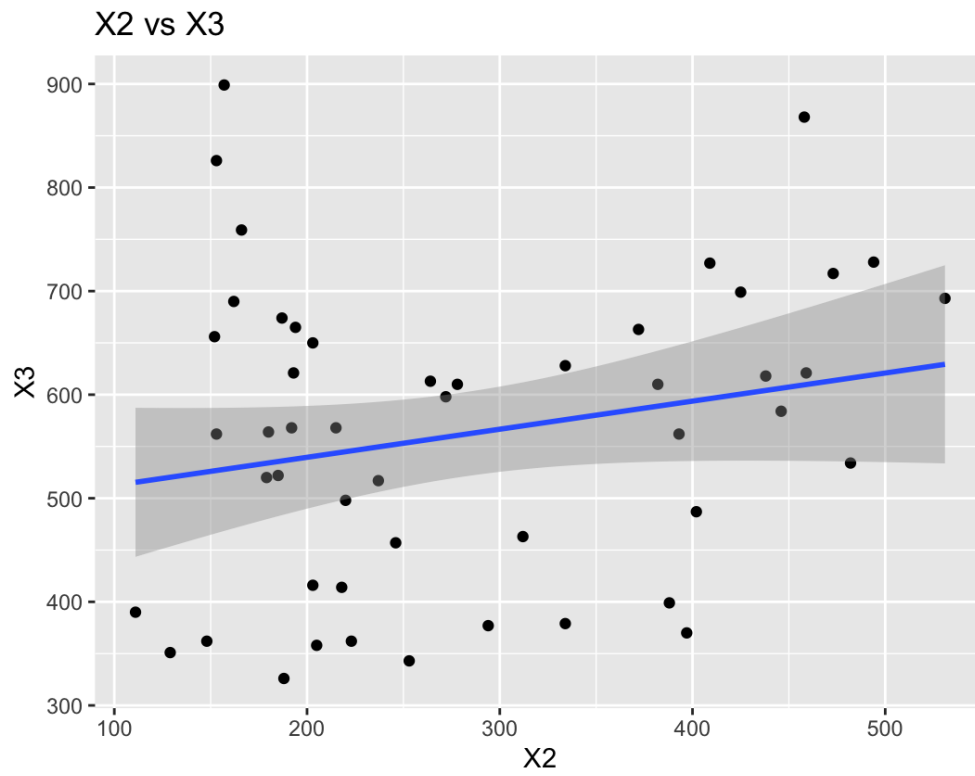


Figure 6: X_2 and X_3 .



Conclusion Based on the correlation matrix and graphs:

1. Y and X1 are positively related. When X1 goes up, Y tends to also go up.
2. Y and X2 have a moderate positive relationship. When X2 increases, Y generally goes up, but not as consistently as with X1.
3. Y and X3 also have a moderate positive relationship. Similar to X2, when X3 rises, Y often goes up, but not all the time.
4. X1 and X2 have a low positive relationship. They somewhat move in the same direction, but it's not very consistent.
5. X1 and X3 are positively related. When one increases, the other usually does too.
6. X2 and X3 have a low positive relationship. They somewhat move in the same direction, but it's not very consistent.

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

Step 1: Calculate the average per capita expenditure on housing assistance for each region

```
1 mean_expenditure_by_region <- aggregate(Y ~ Region, data =  
  expenditure, FUN = mean)  
2  
3 print(mean_expenditure_by_region)  
4  
5 highest_region <- mean_expenditure_by_region[which.max(mean_  
  expenditure_by_region$Y), ]  
6 cat("The region with the highest per capita expenditure on housing  
  assistance is:", highest_region$Region)
```

	Region	Y
1	1	79.44444
2	2	83.91667
3	3	69.18750
4	4	88.30769

The region with the highest per capita expenditure on housing assistance is: 4

Step 2: Plot the relationship between Y and Region

```

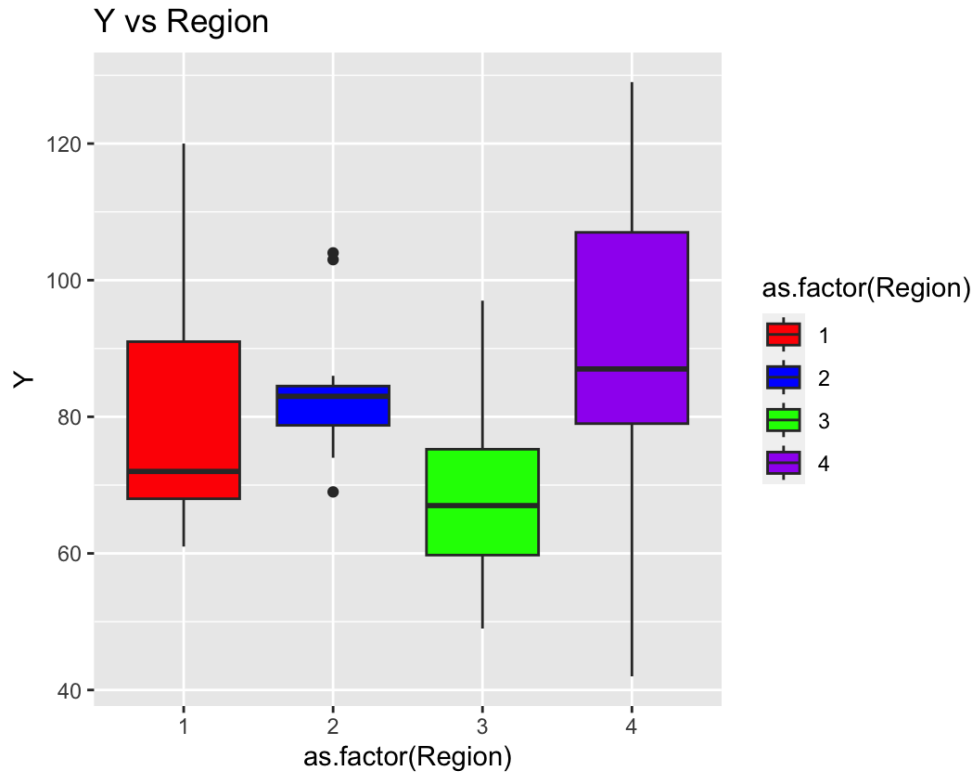
1 highest_region <- mean_expenditure_by_region[which.max(mean_
  expenditure_by_region$Y), ]
2 cat("The region with the highest per capita expenditure on housing
  assistance is:", highest_region$Region)
3
4 plot7<- ggplot(expenditure, aes(x = as.factor(Region), y = Y, fill
  = as.factor(Region))) +
5   geom_boxplot() +
6   ggtitle("Y vs Region") +
7   scale_fill_manual(values = c("red", "blue", "green", "purple"))
8 print(plot7)

```

Description

1. Northeast Region (Red): This region has relatively stable and moderately high expenditure on average.
2. North Central Region (Blue): Expenditure levels are moderate and relatively consistent in this region.
3. South Region (Green): Expenditure levels vary significantly, with some areas having low expenditure and others having high expenditure.
4. West Region (Purple): On average, this region has the highest expenditure, but there is also some variation in expenditure levels.

Figure 7: Y and Regions.



In summary, the West region(4) has the highest average expenditure, the South region(3) shows the most variation, and the Northeast(1) and North Central regions(2) have moderate expenditure levels.

3.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.

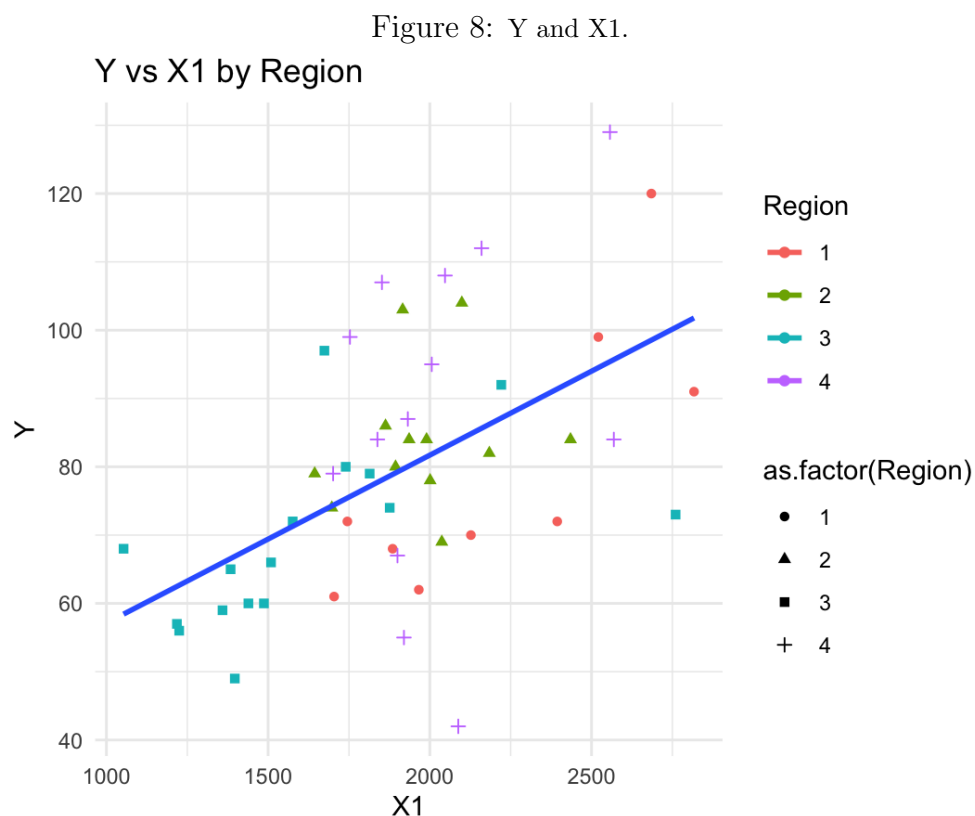
Plot the relationship between Y, as shown in Figure 1.

Reproduce the graph of Y and X1.

```

1 plot8 <- ggplot(expenditure, aes(x = X1, y = Y, color = as.factor(
  Region), shape = as.factor(Region))) +
2   geom_point() +
3   geom_smooth(method = 'lm', se = FALSE, aes(group = 1)) +
4   ggtitle("Y vs X1 by Region") +
5   xlab("X1") +
6   ylab("Y") +
7   labs(color = "Region") +
8   theme_minimal()
9
10 print(plot8)

```



From the graph, it appears that there is a linear relationship between Y and X1. Therefore, I proceed to conduct a statistical test in R to confirm this.

```

1 save_lm_summary <- function(model, file_name) {
2   summary_text <- capture.output(summary(model))
3   writeLines(summary_text, con = file_name)
4 }
5 # Run the linear regression model between Y and X1
6 linear_regression_Y_X1 <- lm(Y ~ X1, data = expenditure)
7
8 linear_regression_summary <- summary(linear_regression_Y_X1)
9
10 linear_regression_summary_text <- capture.output(linear_regression_
    summary)
11 # Save the summary to a .tex file
12 writeLines(linear_regression_summary_text, "/Users/miadong/Desktop/
    linear_regression_Y_X1_summary.tex")
13 # Load the stargazer package
14 library(stargazer)
15
16 # Run the regression
17 regression1 <- lm(Y ~ X1, data = expenditure)
18
19 # Define the output_stargazer function
20 output_stargazer <- function(outputFile, ...) {
21   output <- capture.output(stargazer(..., type = "latex"))
22   cat(paste(output, collapse = "\n"), "\n", file = outputFile,
        append = TRUE)
23 }
24
25 # Use the function to save the LaTeX-formatted table to a .tex file
26 output_stargazer("regression_output11.tex", regression1)
27 getwd()

```

Table 1:

<i>Dependent variable:</i>	
	Y
X1	0.025*** (0.006)
Constant	32.546*** (11.034)
Observations	50
R ²	0.283
Adjusted R ²	0.268
Residual Std. Error	15.836 (df = 48)
F Statistic	18.920*** (df = 1; 48)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Based on the data presented in the table, a linear regression model is evident between Y and $X1$, represented by the equation $Y = 32.546 + 0.025 \cdot X1$.