

# Problem Set 1

## Applied Stats/Quant Methods 1

Due: October 1, 2023

### 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 Sunday October 1, 2023. No late assignments will be accepted.
- Total available points for this homework is 80.

### Question 1 (40 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. Find a 90% confidence interval for the average student IQ in the school.

```
1 # Calculation using standard normal distribution, but the sample size was small
2 # confidence_level <- 0.90
3 # z90 <- qnorm((1 - confidence_level) / 2, lower.tail = FALSE)
4 # sample_data <- y
5 # n <- length(sample_data)
6 # sample_mean <- mean(sample_data)
7 # sample_sd <- sd(sample_data)
8 # lower_90 <- sample_mean - (z90 * (sample_sd / sqrt(n)))
9 # upper_90 <- sample_mean + (z90 * (sample_sd / sqrt(n)))
10 # confint90 <- c(lower_90, upper_90)
11 # cat("90% Confidence Interval:", lower_90, "-", upper_90, "\n")
12 # Calculate it with the t distribution
13 sample_mean <- mean(y)
14 sample_sd <- sd(y)
15 sample_size <- length(y)
16 confidence_level <- 0.90
17 degrees_of_freedom <- sample_size - 1
18 t_critical <- qt(1 - (1 - confidence_level) / 2, df = degrees_of_freedom)
19 standard_error <- sample_sd / sqrt(sample_size)
20 confidence_interval <- c(sample_mean - t_critical * standard_error, sample_mean + t_
    critical * standard_error)
21 cat("90% Confidence Interval for Average IQ:", confidence_interval, "\n")
```

0.90 Confidence Interval for Average IQ: 93.95993 102.9201

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.

```
1 # Hypothesis testing:
2 # Null hypothesis (H0) : The average student IQ of the school is less than or equal to
  the national average IQ score.
3 # Alternative hypothesis (H1) : The average student IQ in the school is greater than
  the national average IQ score.
4 # The alternative hypothesis (H1) describes my research hypothesis ,
5 # so it focuses on a 'greater than' scenario.
6 # Therefore , a right/upper-tailed hypothesis test should be used ,
7 # with p-value calculated as p_value = 1 - pt(t_statistic , df = degrees_of_freedom)."
8 # Make a decision
9 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)
10 sample_mean <- mean(y)
11 sample_sd <- sd(y)
12 sample_size <- length(y)
13 confidence_level <- 0.05
14 degrees_of_freedom <- sample_size - 1
15 t_statistic <- (sample_mean - 100) / (sample_sd / sqrt(sample_size))
16 p_value <- 1 - pt(t_statistic , df = degrees_of_freedom)
17 if (p_value <= confidence_level) {
18   cat("Reject the null hypothesis. The average student IQ is higher than the national
     average IQ score")
19 } else {
20   cat("Fail to reject the null hypothesis. The average student IQ is less than or equal
     to the national average IQ score")
21 }
```

Fail to reject the null hypothesis. The average student IQ is less than or equal to the national average IQ score

## Question 2 (40 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.

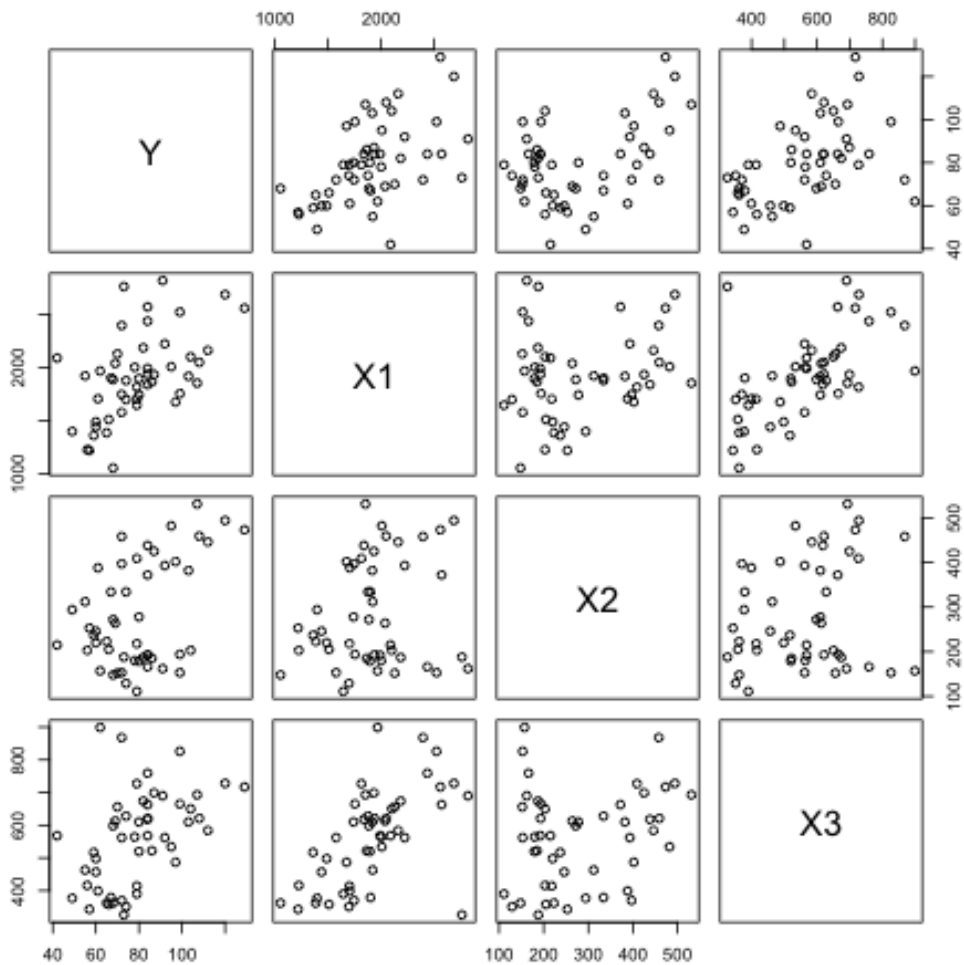
1. Please plot the relationships among  $Y$ ,  $X1$ ,  $X2$ , and  $X3$ ? What are the correlations among them (you just need to describe the graph and the relationships among them)?

```
1 # read in expenditure data
2 expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI_Fall2023/main/datasets/expenditure.txt", header=T)
3
4 # 1
5 data <- read.table(text = "
6 STATE Y X1 X2 X3 Region
7 ME 61 1704 388 399 1
8 NH 68 1885 272 598 1
9 VT 72 1745 397 370 1
10 MA 72 2394 458 868 1
11 RI 62 1966 157 899 1
12 CT 91 2817 162 690 1
13 NY 120 2685 494 728 1
14 NJ 99 2521 153 826 1
15 PA 70 2127 152 656 1
16 OH 82 2184 187 674 2
17 IN 84 1990 192 568 2
18 IL 84 2435 166 759 2
19 MI 104 2099 203 650 2
20 WI 84 1936 193 621 2
21 MN 103 1916 382 610 2
22 IA 86 1863 185 522 2
23 MO 69 2037 264 613 2
24 ND 74 1697 129 351 2
25 SD 79 1644 111 390 2
26 NB 80 1894 179 520 2
27 KS 78 2001 180 564 2
28 DE 73 2760 188 326 3
29 MD 92 2221 393 562 3
30 VA 97 1674 402 487 3
31 WV 66 1509 205 358 3
32 NC 65 1384 223 362 3
33 SC 57 1218 253 343 3
34 GA 60 1487 220 498 3
35 FL 74 1876 334 628 3
36 KY 49 1397 294 377 3
37 TN 60 1439 246 457 3
38 AL 59 1359 237 517 3
39 MS 68 1053 148 362 3
40 AR 56 1225 203 416 3
```

```

41 LA 72 1576 153 562 3
42 OK 80 1740 278 610 3
43 TX 79 1814 409 727 3
44 MT 55 1920 312 463 4
45 ID 79 1701 218 414 4
46 WY 42 2088 215 568 4
47 CO 108 2047 459 621 4
48 NM 84 1838 438 618 4
49 AZ 87 1932 425 699 4
50 UT 99 1753 194 665 4
51 NV 84 2569 372 663 4
52 WA 112 2160 446 584 4
53 OR 95 2006 482 534 4
54 CA 129 2557 473 717 4
55 AK 67 1900 334 379 4
56 HI 107 1852 531 693 4
57 ", header = TRUE)
58 install.packages("ggplot2")
59 plot(expenditure[c("Y", "X1", "X2", "X3")])
60
61 # 0.5317212 Y and X1: Positively correlated, strong relationship
62 cor(expenditure["Y"], expenditure["X1"])
63
64 # 0.4482876 Y and X2: Positively correlated, moderate strength of relationship
65 cor(expenditure["Y"], expenditure["X2"])
66
67 # 0.4636787 Y and X3: Positively correlated, moderate strength of relationship.
68 cor(expenditure["Y"], expenditure["X3"])
69
70 # 0.2056101 X1 and X2: Positively correlated, weak relationship.
71 cor(expenditure["X1"], expenditure["X2"])
72
73 # 0.5952504 X1 and X3: Positively correlated, strong relationship.
74 cor(expenditure["X1"], expenditure["X3"])
75
76 # 0.2210149 X2 and X3: Positively correlated, weak relationship.
77 cor(expenditure["X2"], expenditure["X3"])

```



Y and X1: Positively correlated, strong relationship.

Y and X2: Positively correlated, moderate strength of relationship.

Y and X3: Positively correlated, moderate strength of relationship.

X1 and X2: Positively correlated, weak relationship.

X1 and X3: Positively correlated, strong relationship.

X2 and X3: Positively correlated, weak relationship.

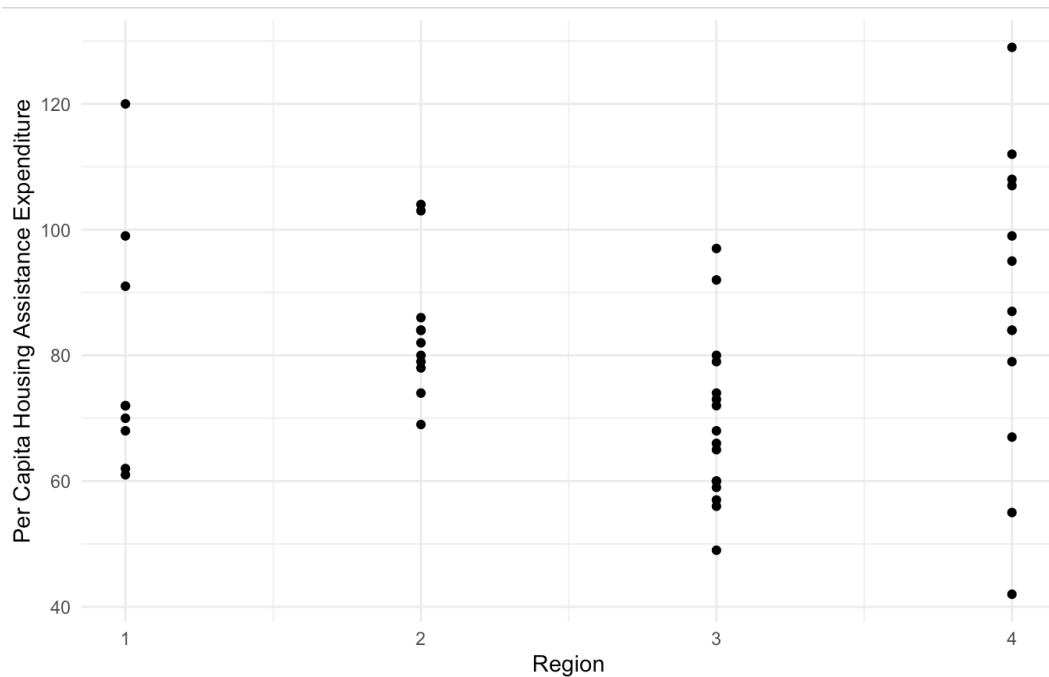
2. Please plot the relationship between Y and Region? On average, which region has the highest per capita expenditure on housing assistance?

```
1 # 2
2 library(ggplot2)
3 # Create a scatter plot to show the relationship between Y and Region
4 # Calculate the average Y value by Region
```

```

5 # Find the region with the highest average Y value
6 ggplot(data, aes(x = Region, y = Y)) +
7   geom_point() +
8   labs(x = "Region", y = "Average Housing Assistance Expenditure (Y)") +
9   theme_minimal()
10
11 average_by_region <- aggregate(data$Y, by = list(data$Region), FUN = mean)
12 max_region <- average_by_region[which.max(average_by_region$x), ]
13 cat("On average, the areas with the highest per capita expenditure on housing assistance are
: ", max_region$Group.1, "\n")

```



On average, the areas with the highest per capita expenditure on housing assistance are: 4

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.

```

1 # 3
2 # Plot a scatter plot of the relationship between Y and X1,
3 # using different types of symbols and colors for different regions
4 ggplot(data, aes(x = X1, y = Y, shape = as.factor(Region), color = as.factor(Region))) +
5   geom_point() +
6   labs(x = "per capita personal income in state", y = "housing assistance in state") +
7   theme_minimal()
8 # 0.5317212 Y and X1: Positively correlated, strong relationship
9 cor(expenditure["Y"], expenditure["X1"])

```

Y and X1: Positively correlated, strong relationshi

