# Problem Set 1

Valeriia Babaian

Due: October 1, 2023

## Question 1 (40 points): Education

1. Find a 90% confidence interval for the average student IQ in the school.
   *Answer*: [93.95993; 102.92007].

```
1  t_score <- qt(0.95, df=length(y)-1)
2  lower_90_t <- mean(y)-(t_score)*(sd(y)/sqrt(length(y)))
3  upper_90_t <- mean(y)+(t_score)*(sd(y)/sqrt(length(y)))
4
5  # Confidence interval boundaries: 93.95993; 102.9201
6  lower_90_t
7  mean(y) #98.44
8  upper_90_t
9
10 #double checking:
11 t.test(y, conf.level = 0.9, alternative = "two.sided")
12 # 90 percent confidence interval: 93.95993; 102.92007
```

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.
   *Answer*: It is not higher:

```
1  t.test(y, mu = 100, alternative = "greater")
2  # We cannot reject the null hypothesis that the average student IQ in the
       school is not higher than the average IQ score across all schools in
       the country (p-value = 0.7215)
```

Using the same sample, conduct the appropriate hypothesis test with $\alpha = 0.05$.
*Answer*: Expectedly, the answer is the same as the average student IQ score in the given school was not greater than average IQ across the country on any of conventional levels of confidence.

```
1 ##Using the same sample, conduct the appropriate hypothesis test with \
       alpha = 0.05.
2 t.test(y, mu = 100, conf.level = 0.95) # we cannot first reject the H0
       that the averages are equal on 0.05 significance (p-value = 0.5569)
```

# Question 2 (40 points): Political Economy

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

```
1 expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI_
    Fall2023/main/datasets/expenditure.txt", header=T)
2 str(expenditure)
```

- 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 library(GGally)
2 library(ggplot2)
3 ggpairs(expenditure,
4        columns = 2:5,
5        columnLabels=c("Housing assistance in state",
6     "Personal income in state",
7     "Financially insecure share",
8     "Urban population share"))
9 ggsave("Babaian-plot1.pdf", width = 8,   height = 6)
```
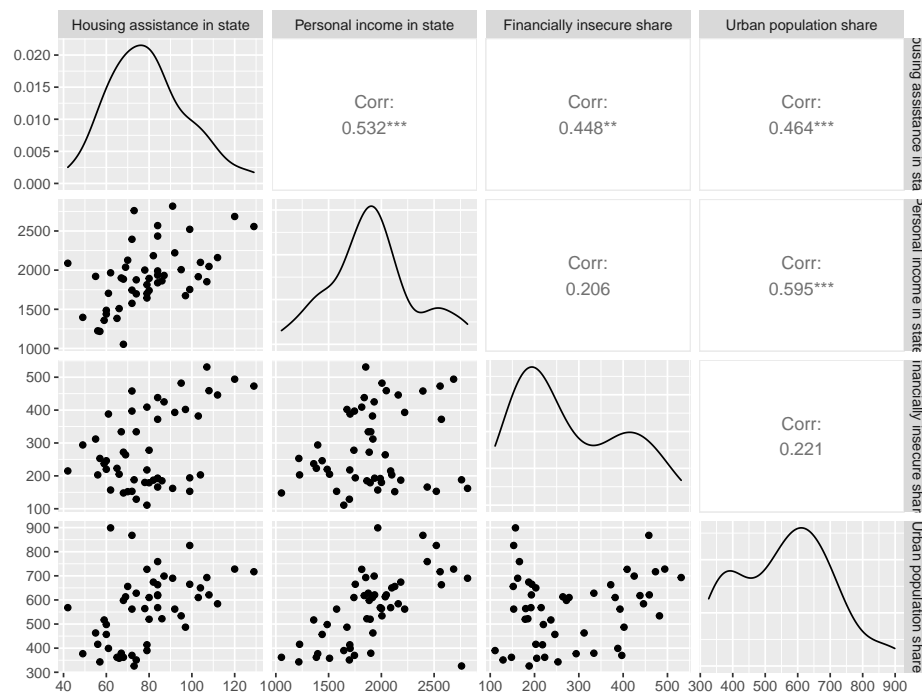


Figure 1: Fig. 1

According to the visualized relationships between all meaningful quantities in the data, the variables are positively related, but there are no strong correlations between any of them. The upper-right part of the graphics proves this point.

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

```
1  expenditure$Region <- factor(expenditure$Region, labels = c("Northeast",
     "North Central", "South", "West"))
2
3  ggplot(data=expenditure, mapping=aes(x=Region, y=Y)) +
4    stat_summary(fun.data=mean_sdl, geom="bar")   +
5    labs(title="Per capita expenditure on shelters/housing assistance in
     state by region",
6         x="", y="")
7  ggsave("Babaian-plot2.pdf", width = 8,    height = 6)
```
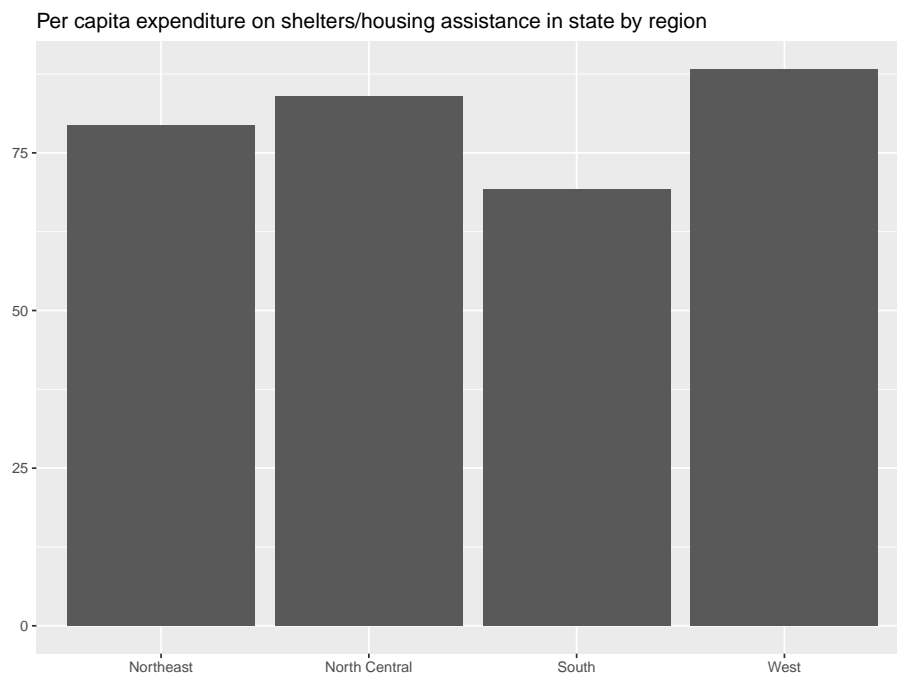


Figure 2: Fig. 2

The highest per capita expenditure on housing assistance is in the Western states.

- 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

4

different regions with different types of symbols and colors.

```
1  plot(expenditure$X1, expenditure$Y,
2      xlab="Per capita personal income in state",
3      ylab="Housing assistance in state",
4      main="The Relationship between state-level personal income and
   housing assistance")
5  ggsave("Babaian-plot3.pdf", width = 8,   height = 6)
6  # There is a positive relationship between per capita personal income in
   state
```
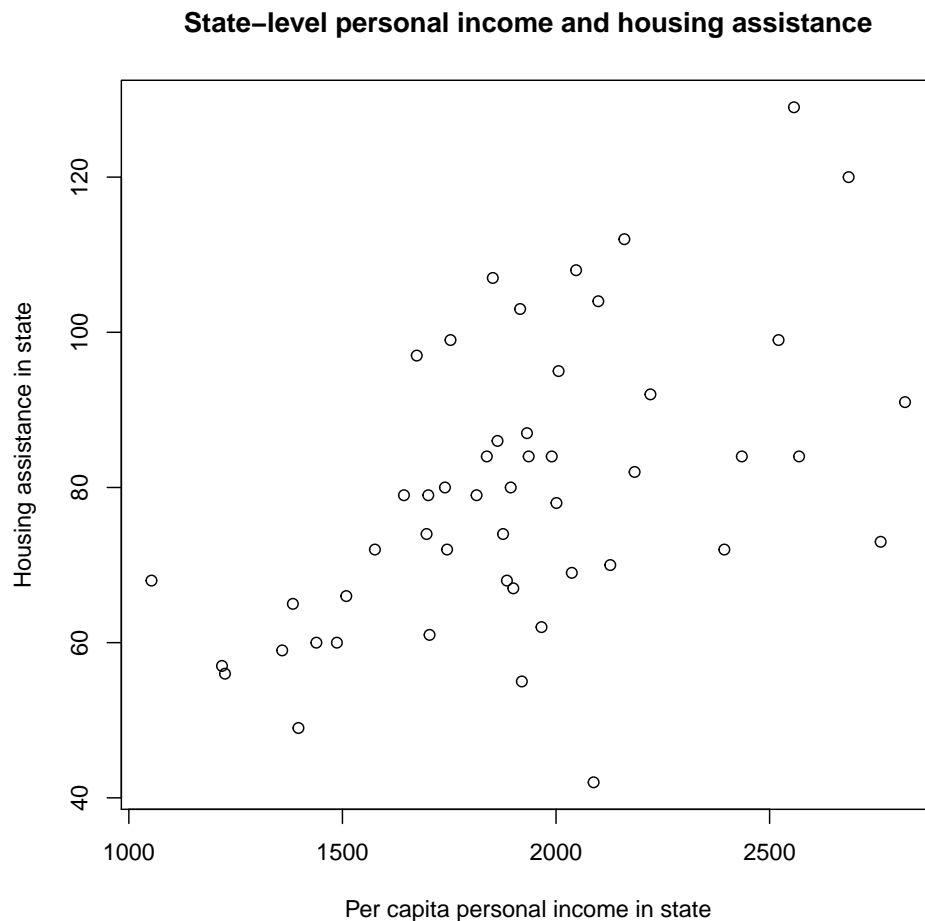


Figure 3: Fig. 3

There is a positive relationship between per capita personal income in state and its housing assistance expenditures, however, there is no strong correlation as there is a significant amount on deviations on both sides of the distributions.

```
1  plot(expenditure$X1, expenditure$Y,
2      col= expenditure$Region,
3      pch=c(21,22,23,24)[as.numeric(expenditure$Region)],
4      xlab="Per capita personal income in state",
5      ylab="Housing assistance in state",
6      main="State-level personal income and housing assistance by regions")
7  legend("topleft", cex=1, legend=levels(expenditure$Region),
8          pch=unique(expenditure$Region), col=unique(expenditure$Region))
9  dev.off()
```



Figure 4: Fig. 4