

Problem Set 1

Imelda Finn, 22334657

Due: October 3, 2021

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 iqData <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113,  
2           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.

```
# calculate sample statistics  
# capture the number of observations  
n <- length(iqData)  
  
# calculate mean  
iqSum <- sum(iqData)           # sum of IQ scores  
iqMean <- iqSum / n           # mean IQ score for sample  
  
# calculate variance and standard deviation  
iqVar <- sum((iqData - iqMean)^2)/(n-1)  
iqSD <- sqrt(iqVar)  
  
iqse <- iqSD / sqrt(n)         # standard error of sample
```

The code for the t-test at 90% is:¹

```
t.val <- qt(alphaVal/2, df = n-1, lower.tail = FALSE)

CI_lower <- iqMean - t.val * iqse
CI_upper <- iqMean + t.val * iqse
```

The result was:

Our Confidence interval for the IQ of the students in the sample is:

93.96 < mean IQ < 102.92

with a confidence level of 90%. ²

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 $\alpha = 0.05$.*

- (a) The number of observations is 25, which isn't ideal for t-test statistics (we would prefer at least 30 observations). The data isn't really random, as all are students from the same school.
- (b) H_0 : the average IQ score in the sample is less than the population average ie $\mu_O \leq \mu$
- (c) H_a : the average IQ score in the sample is greater than the population average ie $\mu_O > \mu$
- (d) Calculate test statistic

$$TS = \frac{\bar{Y} - \mu_0}{\sigma_{\bar{Y}}}$$

- (e) Calculate p-value

$$p = Pr(Z \leq -|\frac{\bar{Y} - \mu_O}{\sigma_{\bar{Y}}}|)$$

- (f) if $p \leq \alpha = 0.05$, we reject the null hypothesis

```
# Test our hypothesis
alphaVal <- 0.05
```

¹t.val = 1.71

²A Z-test gave a 90% CI of $94.13 < \mu < 102.75$.

```

popMean <- 100
#get test statistic
testStatistic <- (iqMean - popMean) / iqse

pValue <- pnorm(-abs(testStatistic))

# calculate t-test p-value
t_pValue <- pt(abs(testStatistic), df = n-1, lower.tail = FALSE)

```

Results

p-value for normal distribution is 0.276

t	df	p-value	
-0.5957439	24	0.2784617	The p-value is greater than 5%, so we cannot reject the null hypothesis. The data does not support the suggestion that the school IQ scores are greater than the population average.

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

```
1 #expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI-Fall2022/main/datasets/expenditure.txt", header=T)
```

- In Figure 1, there appears to be some positive correlation between the per capita expenditure on Shelters/Housing Assistance (SHA) and each of the three numerical variables. In each case, an increase in the response variable accompanies an increase in the candidate explanatory variable.

The weakest relationship is between SHA expenditure and the number of residents who are financially insecure. When looked at in more detail, Figure 2 shows that there is a negative correlation between the two variables (ie SHA expenditure per capita decreases as the number of financially insecure residents increases) until $X2$ reaches approximately 300 per 100,000; after that point, the values increases are positively correlated. (line fitted using: 'geom_smooth()' using method = 'loess' and formula 'y ~ x')

The income per capita shows some correlation with expenditure and with the number of urban residents, but not with financially insecure residents. Similarly, the plot of urban residents against financially insecure residents appears to be uniformly distributed, (i.e. uncorrelated), while expenditure and income both increase with the increase in urban population. Figure 3 shows that the slope of the line relating expenditure on SHA to the urban population turns negative after $X3$ reaches approximately 750 per 1,000, however there are very few data points in this region of the plot.

-

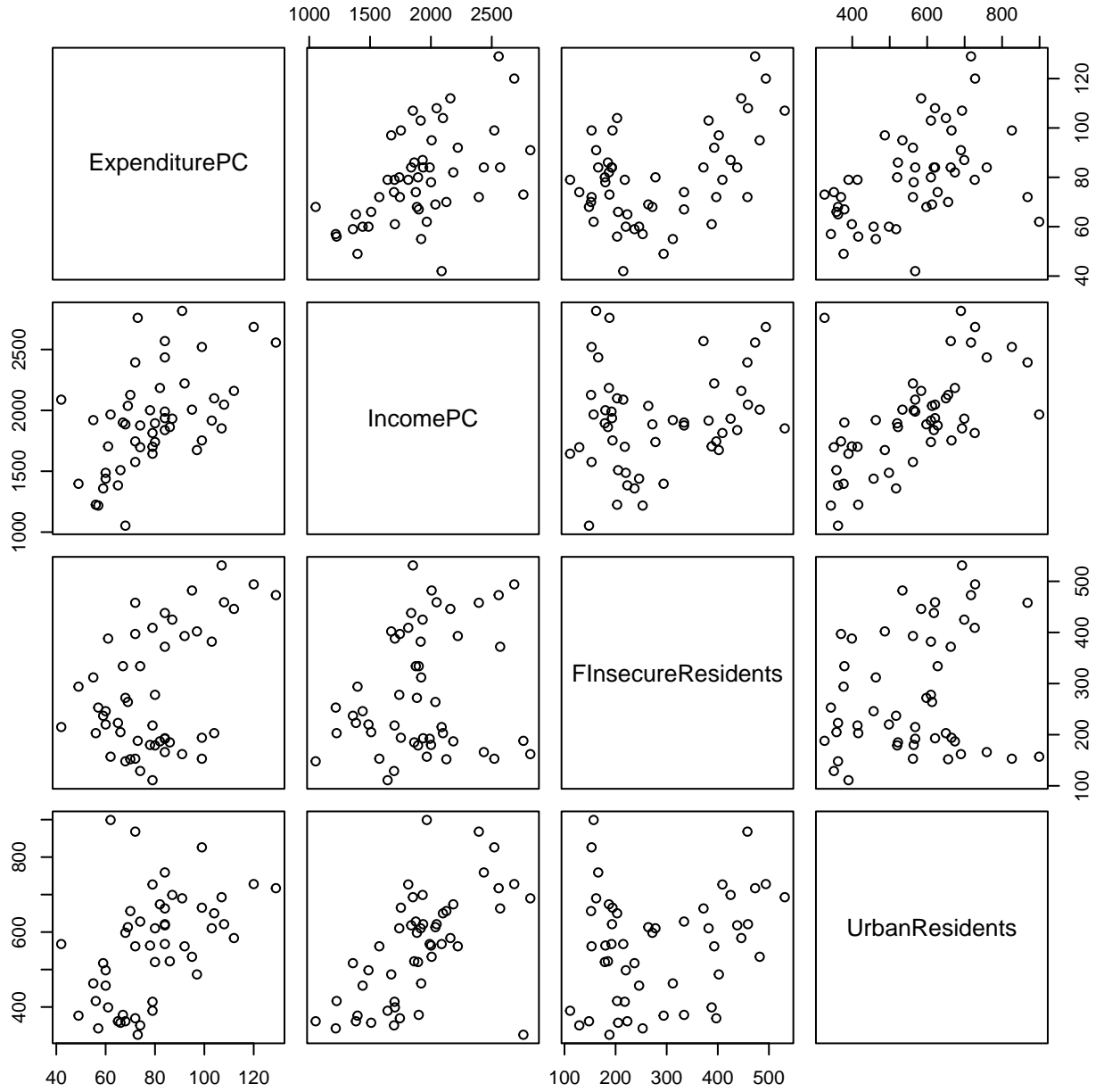


Figure 1: Pair-wise comparison of expenditure variables

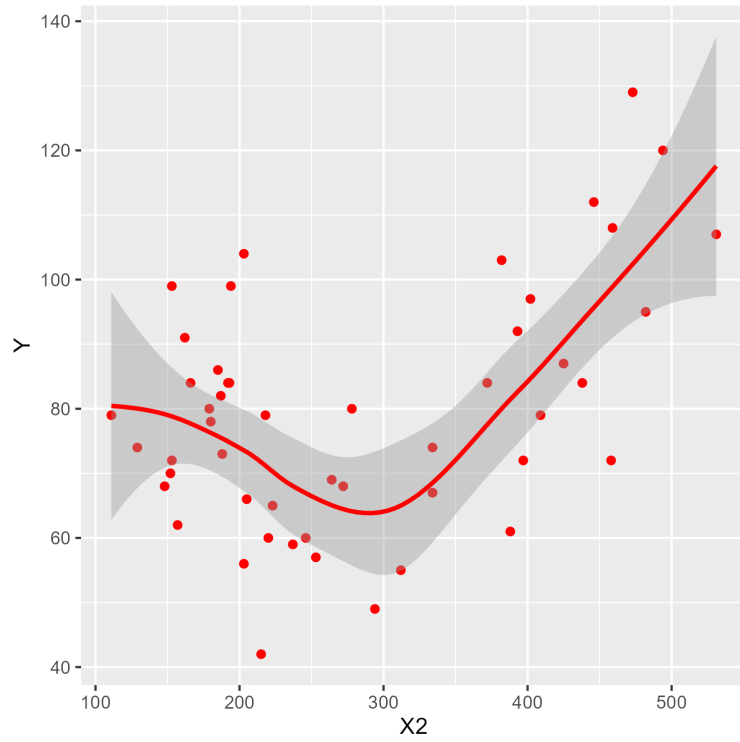


Figure 2: Relationship between expenditure on SHA and financially insecure population

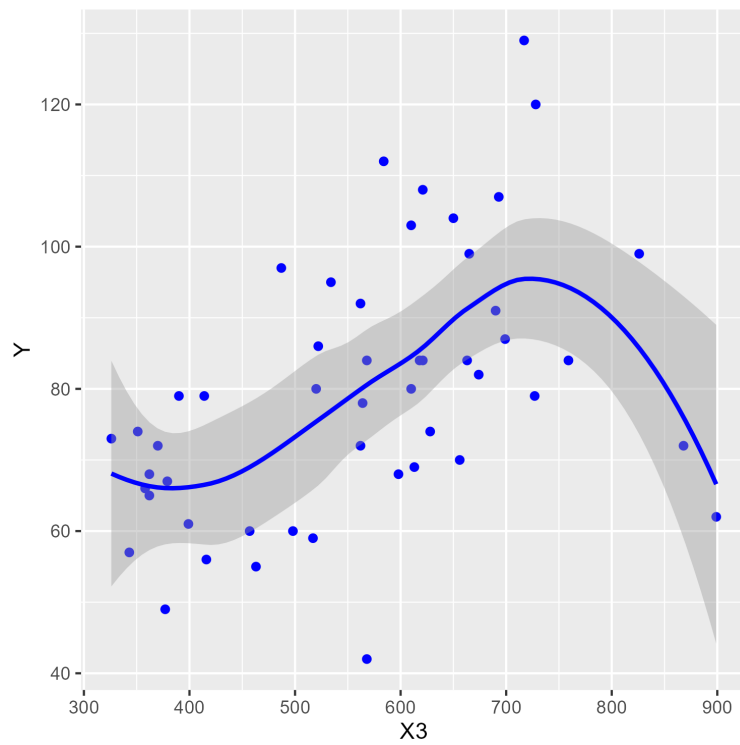


Figure 3: Relationship between expenditure on SHA and urban population

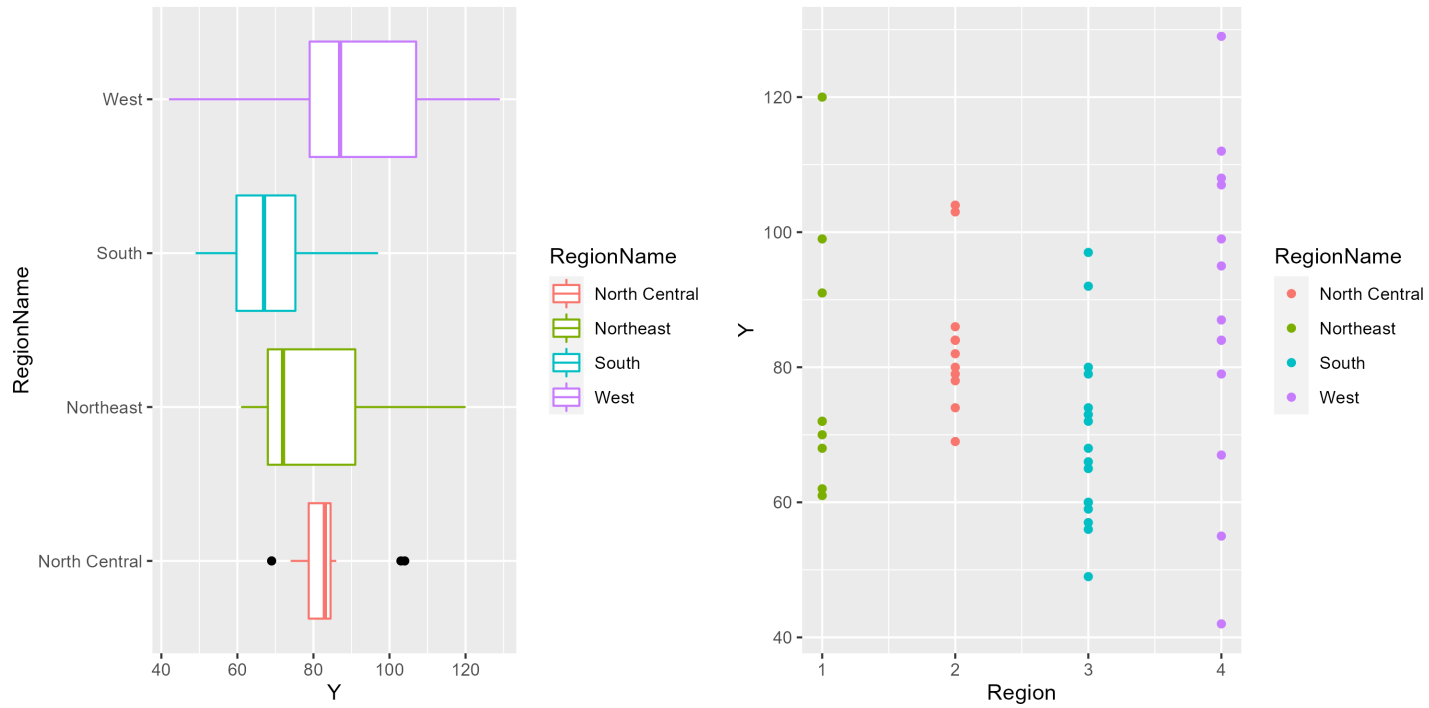


Figure 4: SHA expenditure by geographical region

The relationship between expenditure on SHA and region is shown in Figure 4

```
regional_mean_table <-expenditure %>% # Tidyverse method for grouping
  group_by(Region) %>%
  summarise(mean = round(mean(Y), 2))
```

Table 1: Regional spending on SHA

	Region	mean	regions
1	1	79.44	Northeast
2	2	83.92	North Central
3	3	69.19	South
4	4	88.31	West

The West region has the highest average per capita spending (\$88.31) on Shelters/Housing Assistance (Table 1).

- The graph of income per capita (X1, \$) against expenditure on SHA per capita (Y, \$) appears to show a positive correlation between the two variables, i.e. it suggests that higher pc income could be a predictor for higher spending on housing assistance (Figure 5).

When the data is subdivided by region the relationship between the variables is much less convincing. The individual regions show clear differences between the interaction of the two values - North Central and South show no positive correlation; in the Northeast and West it is possible that some kind of relationship exists, but it's not consistent (Figure 6).

The apparent relationship when the data are combined clearly masks underlying differences between the regions.

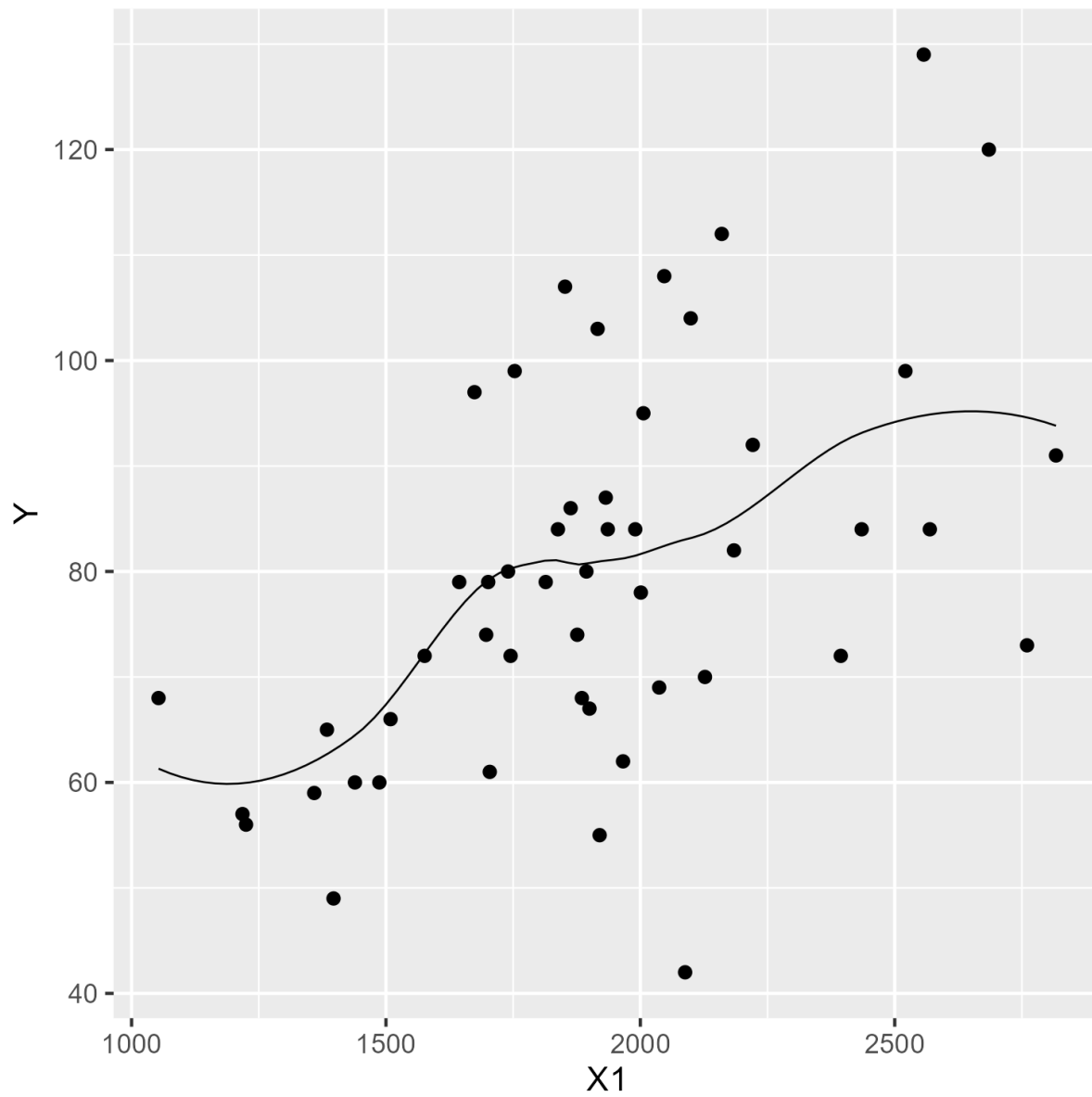


Figure 5: Income vs SHA expenditure, per capita

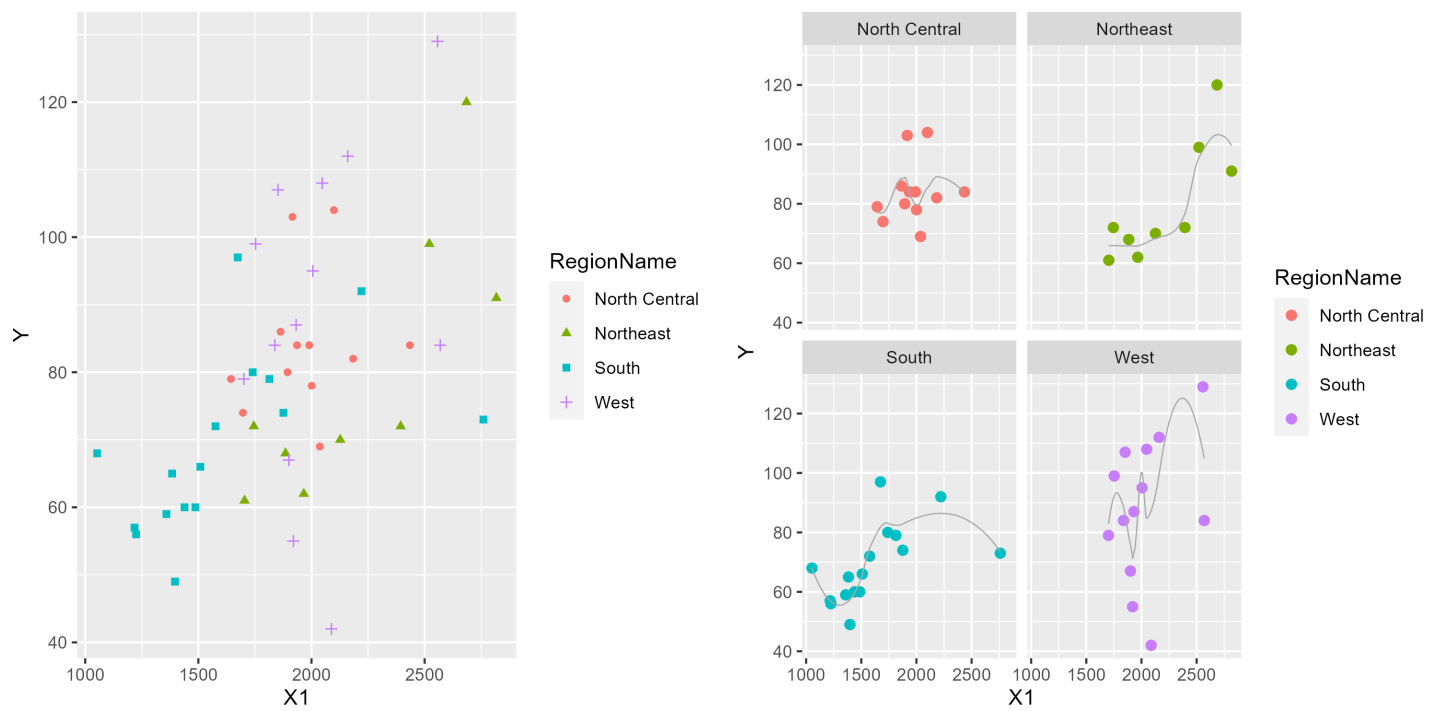


Figure 6: Income per capita vs SHA expenditure, by Region

Appendix - R code

```
1 #####
2 # Imelda Finn, 22334657
3 # POP77003 – Stats I
4 # clear global .envir, load libraries, set wd
5 #####
6
7 # remove objects
8 rm(list=ls())
9
10 # detach all libraries
11 detachAllPackages <- function() {
12   basic.packages <- c("package:stats", "package:graphics", "package:grDevices"
13     , "package:utils", "package:datasets", "package:methods", "package:base")
14   package.list <- search()[ifelse(unlist(gregexpr("package:", search()))==1,
15     TRUE, FALSE)]
16   package.list <- setdiff(package.list, basic.packages)
17   if (length(package.list)>0) for (package in package.list) detach(package,
18     character.only=TRUE)
19 }
20 detachAllPackages()
21
22 # load libraries
23 pkgTest <- function(pkg){
24   new.pkg <- pkg[!(pkg %in% installed.packages()[, "Package"])]
25   if (length(new.pkg))
26     install.packages(new.pkg, dependencies = TRUE)
27   supply(pkg, require, character.only = TRUE)
28 }
29
30 # load necessary packages
31 lapply(c("ggplot2", "stargazer", "tidyverse", "stringr", "quantreg"), pkgTest
32 )
33
34 # set working directory to current parent folder
35 setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
36
37 #####
38 # Problem 1
39 #####
40
41 # load data as vector – in .tex file – update if move from 38
42 iqData <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113,
43   112, 98, 80, 97, 95, 111, 114, 89, 95, 126, 98)
44
45 ## Save our data to a .csv file in the data directory
46 write.csv(iqData, file = "Data/iq.csv", row.names = FALSE)
47
48 # Explore data
49 summary(iqData)
50 str(iqData)
```

```

47 head(iqData)
48
49 # look at sampling from sample
50 meanIQ <- vector("double", length = 1000)
51 for (i in 1:1000) {
52   meanIQ[i] <- mean(sample(iqData, 25, replace=TRUE))
53 }
54 summary(meanIQ)
55 boxplot(meanIQ, iqData, xlab=c("averaged vs original sample"))
56
57
58 # Visually inspect the data
59 hist(iqData, breaks = 10, main = "Histogram of IQ", xlab = "IQ")
60
61 plot(density(iqData), main = "PDF of IQ", xlab = "IQ")
62
63 # Use a QQ plot to determine if our IQ variable is normally distributed
64 qqnorm(iqData)
65 qqline(iqData, distribution = qnorm)
66 # Sample values fall away from normal line at upper end
67
68 ##-----
69 # calculate sample statistics
70 # capture the number of observations
71 n <- length(iqData)
72
73 # calculate mean
74 iqSum <- sum(iqData)           # sum of IQ scores
75 iqMean <- iqSum / n           # mean IQ score for sample
76
77 # calculate variance and standard deviation
78 iqVar <- sum((iqData - iqMean)^2)/(n-1)
79 iqSD <- sqrt(iqVar)
80
81 iqse <- iqSD / sqrt(n)         # standard error of sample
82
83 ##-----
84 ## Confidence Intervals
85 # Calculate 90 percent confidence intervals using normal distribution
86 # assuming iqMean ~ N(mu, iqse)
87 alphaVal = 0.1
88 CI_lower <- qnorm(alphaVal/2, mean = iqMean, sd = iqse)
89
90 CI_upper <- qnorm(1-alphaVal/2, mean = iqMean, sd = iqse)
91
92 # output
93 cat(str_glue("{(1-alphaVal)*100}% Confidence Intervals, two-sided z-test"))
94 matrix(c(CI_lower, CI_upper), ncol = 2,
95        dimnames = list("", c("Lower", "Upper")))
96
97 # Calculate 90 percent confidence intervals using t-distribution

```

```

98 # degrees of freedom = n-1 = 24 - should be >30
99 t.val <- qt(alphaVal/2, df = n-1, lower.tail = FALSE)
100
101 CI_lower <- iqMean - (t.val * iqse)
102 CI_upper <- iqMean + (t.val * iqse)
103
104 # calculate using t-test
105 cat(str_glue("{(1-alphaVal)*100}% Confidence Intervals , two-sided t-test"))
106 matrix(c(CI_lower, CI_upper), ncol = 2,
107        dimnames = list("", c("Lower", "Upper")))
108
109 # t-test results in (slightly) wider confidence interval
110
111 # Check our working
112 #t.test(iqData, conf.level = 1-alphaVal, alternative = "two.sided")
113
114 cat("Our Confidence interval for the IQ of the students in the sample is: ")
115 cat(str_glue(" {round(CI_lower,2)} < mean IQ < {round(CI_upper,2)} "))
116 cat(str_glue("with a confidence level of {(1-alphaVal)*100}%"))
117
118
119 ##-----
120 ## Hypothesis Testing
121 # Wrangling our data
122 class(iqData) # What class of vector is our IQ variable? - numeric
123
124 # Hypothesis test:
125 # H0 : average IQ of students in school is less than or equal to national
    average
126 # Ha : average IQ of students in school is greater than national average
127
128 # alpha = 0.05, 1-tail test, single population
129
130 # don't have variance of population, only have mean to compare against
131
132 # Test our hypothesis
133 alphaVal <- 0.05
134 popMean <- 100
135 #get test statistic
136 testStatistic <- (iqMean - popMean) / iqse
137
138 pValue <- pnorm(-abs(testStatistic))
139
140 # calculate t-test p-value
141 t_pValue <- pt(abs(testStatistic), df = n-1, lower.tail = FALSE)
142
143 matrix(c(testStatistic, n-1, t_pValue), ncol = 3,
144        dimnames = list("", c("t", "df", "p-value")))
145 cat(str_glue("p-value for normal distribution is {round(pValue,3)}"))
146
147 # check result

```

```

148 t.test( iqData ,
149         mu = 100, # population mean
150         var.equal = TRUE, # The default is FALSE – don't have var for popn
151         alternative = "less", # H0: sample mean > population mean
152         conf.level = .95) #
153
154
155 # How do we interpret the output?
156 # for confidence level of 95%, we cannot reject the hypothesis (p-value >
    alpha)
157
158
159 #####
160 # Problem 2
161 #####
162 rm(list=ls())
163 # set working directory to current parent folder
164 setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
165
166 # function to save output to a file that you can read in later to your docs
167 output_stargazer <- function(outputFile, appendVal=TRUE, ...) {
168   output <- capture.output(stargazer(...))
169   cat(paste(output, collapse = "\n"), "\n", file=outputFile, append=appendVal)
170 }
171
172 # read in expenditure data
173 #expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI-
    Fall2022/main/datasets/expenditure.txt", header=T)
174 #write.table(expenditure, "Data/expenditure.txt")
175 expenditure <- read.table("Data/expenditure.txt", header=T)
176
177 # State 50 states in US
178 #Y per capita expenditure on shelters/housing assistance in state
179 #X1 per capita personal income in state
180 #X2 Number of residents per 100,000 that are "financially insecure" in state
181 #X3 Number of people per thousand residing in urban areas in state
182 #Region 1=Northeast, 2= North Central, 3= South, 4=West
183
184 data_headers <- c("State", "ExpenditurePC", "IncomePC", "FInsecureResidents",
185                  "UrbanResidents", "Region")
186 regions <- c("Northeast", "North Central", "South", "West")
187 names(expenditure)
188 expenditure$RegionName<-regions[expenditure$Region]
189
190 # Inspect the data
191 head(expenditure)
192 str(expenditure)
193 summary(expenditure)
194
195 #investigate spending on Housing assistance
196 # Visualise

```

```

197
198 onefile <- FALSE
199 #pdf( file = if(onefile) "expenditure_plots.pdf" else "expenditure_plots%03d.
    pdf")
200
201 hist(expenditure$Y, main = "Histogram of spending on HA ", xlab = "$, per
    capita")
202
203 plot(density(expenditure$Y),
204       main = "PDF of spending on HA ", xlab = "$, per capita")
205
206 qqnorm(expenditure$Y)
207 qqline(expenditure$Y, distribution = qqnorm)
208
209 #dev.off() # close pdf file
210 #-----
211 # plot the numerical variables against each other
212 pdf("expenditure_pairs.pdf" )
213 pairs(~Y + X1 + X2 + X3, expenditure, labels = data_headers[2:5])
214 dev.off() # close output
215
216
217 # look at detail of some relationships
218 ggplot(expenditure) +
219   geom_point(aes(X2, Y), colour = "red" ) +
220   geom_smooth(aes(X2, Y), colour = "red")
221 ggsave("y-x2.png", width = 5, height = 5)
222
223
224 ggplot(expenditure) +
225   geom_point(aes( X3, Y, ), colour = "blue" ) +
226   geom_smooth(aes( X3, Y), colour = "blue")
227 ggsave("y-x3.png", width = 5, height = 5)
228
229 #-----
230 # regional expenditure on housing assistance
231 # look at plots
232 # by state
233 x <- seq(1, length(expenditure$Y))
234 ggplot(expenditure) +
235   geom_point(aes( x, Y , colour = RegionName))
236
237 # grouped by region
238 ggplot(expenditure) +
239   geom_point(aes(Region, Y, colour = RegionName))
240 ggsave("region-y.png", width = 5, height = 5)
241 # more spread in r4, least in r2
242
243 ggplot(expenditure) +
244   geom_boxplot(aes(Y, RegionName, colour=RegionName), outlier.colour = "black"
    )

```



```

245 ggsave("region_boxplot.png", width = 5, height = 5)
246
247
248 # calculate regional means
249 regional_mean_table <- expenditure %>% # Tidyverse method for grouping
250   group_by(Region) %>%
251   summarise(mean = round(mean(Y), 2))
252
253 regional_mean_table <- cbind(regional_mean_table, regions)
254
255 ggplot(regional_mean_table) +
256   geom_point(aes(regions, mean, colour = regions, shape = regions), size=3)
257 ggsave("region_means.png", width = 5, height = 5)
258
259 # West region has highest per capita mean expenditure on housing assistance
260
261 matrix(regional_mean_table$mean, ncol = 4,
262        dimnames = list("", c(regional_mean_table$regions)))
263 cat(str_glue("Highest average pc spending on SHA is ${regional_mean_table[4,
264   2]} in the West region"))
265
266 output_stargazer("regional_means.tex", appendVal = FALSE, regional_mean_table,
267                 title="Regional spending on SHA", #column.labels=regional_
268                 mean_table$regions,
269                 label="tab:region_mean", summary=FALSE, digits = 2
270                 ) # file fragment
271
272 #
273 # look at income vs expenditure on HA, by region
274 #factor(expenditure$Region) <- regions
275 ggplot(expenditure) +
276   geom_point(aes(X1, Y)) +
277   geom_smooth(aes(X1, Y), colour = "black", se=FALSE, size=0.3)
278 ggsave("y_x1.png", width = 5, height = 5)
279
280 ggplot(expenditure) +
281   geom_point(aes(X1, Y, colour= RegionName, shape= RegionName))
282 ggsave("y_x1_region.png", width = 5, height = 5)
283
284 ggplot(expenditure, ) +
285   geom_point(aes(X1, Y, colour=RegionName), size = 2) +
286   facet_wrap(~ RegionName, nrow = 2) +
287   geom_smooth(aes(X1, Y), colour = "darkgrey", size = 0.3, se=FALSE)
288 ggsave("y_x1_region_facet.png", width = 5, height = 5)

```