# 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:

```
iqData <- 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)
```

    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)
```

The code for the t-test at 90% is:

```
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%. [1]

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).

   (b) $H_0$: the average iq score in the sample is less than the population average ie $\mu_O \leq \mu$

   (c) $H\alpha$: the average iq score in the sample is less than or equal to 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

*Code*

```
# Test our hypothesis
   alphaVal <- 0.05
   popMean <- 100
   #get test statistic
   testStatistic <- (iqMean - popMean) / iqse
```

---

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

```
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* |

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)?

- Please plot the relationship between *Y* and *Region*? On average, which region has the highest 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.

4

Appendix - R code

```r
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"
       , "package:utils", "package:datasets", "package:methods", "package:base")
13   package.list <− search ()[ ifelse ( unlist (gregexpr("package:", search()))==1,
       TRUE, FALSE)]
14   package.list <− setdiff (package.list , basic.packages)
15   if (length(package.list )>0)  for (package in package.list) detach(package,
       character.only=TRUE)
16 }
17 detachAllPackages ()
18
19 # load libraries
20 pkgTest <− function (pkg){
21   new.pkg <− pkg[!(pkg %in% installed.packages ()[,  "Package"])]
22   if (length(new.pkg))
23     install.packages(new.pkg,  dependencies = TRUE)
24   sapply(pkg,  require ,  character.only = TRUE)
25 }
26
27 # load necessary packages
28 lapply(c("ggplot2", "stargazer", "tidyverse", "stringr"),  pkgTest)
29
30 # set working directory to current parent folder
31 setwd(dirname( rstudioapi :: getActiveDocumentContext ()$path ))
32
33 #######################
34 # Problem 1
35 #######################
36
37 # load data as vector  − in .tex file − update if move from 38
38 iqData <− c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113,
39             112, 98, 80, 97, 95, 111, 114, 89, 95, 126, 98)
40
41 ## Save our data to a .csv file in the data directory
42 write.csv(iqData ,
43           file = "Data/iq.csv",
44           row.names = FALSE)
45
46 # Explore data
47 summary(iqData )
```

```r
48  str(iqData)
49  head(iqData)
50
51  # look at sampling from sample
52  meanIQ <- vector("double", length = 1000)
53  for (i in 1:1000) {
54    meanIQ[i] <- mean(sample(iqData, 25, replace=TRUE))
55  }
56  summary(meanIQ)
57  boxplot(meanIQ, iqData, xlab=c("averaged vs original sample"))
58
59
60  # Visually inspect the data
61  hist(iqData, breaks = 10, main = "Histogram of IQ", xlab = "IQ")
62
63  plot(density(iqData), main = "PDF of IQ", xlab = "IQ")
64
65  # Use a QQ plot to determine if our IQ variable is normally distributed
66  qqnorm(iqData)
67  qqline(iqData, distribution = qnorm)
68  # Sample values fall away from normal line at upper end
69
70  ##————————————————————————————————————————
71  # calculate sample statistics
72  # capture the number of observations
73  n <- length(iqData)
74
75  # calculate mean
76  iqSum <- sum(iqData)              # sum of IQ scores
77  iqMean <-iqSum / n               # mean IQ score for sample
78
79  # calculate variance and standard deviation
80  iqVar <- sum((iqData - iqMean)^2)/(n-1)
81  iqSD <-sqrt(iqVar)
82
83  iqse <- iqSD / sqrt(n)           # standard error of sample
84
85  ##————————————————————————————————————————
86  ## Confidence Intervals
87  # Calculate 90 percent confidence intervals using normal distribution
88  # assuming  iqMean ~ N(mu, iqse)
89  alphaVal = 0.1
90  CI_lower <- qnorm(alphaVal/2, mean = iqMean, sd =  iqse)
91
92  CI_upper <- qnorm(1-alphaVal/2, mean = iqMean, sd = iqse)
93
94  # output
95  cat(str_glue("{(1-alphaVal)*100}% Confidence Intervals, two-sided z-test"))
96  matrix(c(CI_lower, CI_upper), ncol = 2,
97         dimnames = list("",c("Lower", "Upper")))
98
```

```r
 99 # Calculate 90 percent confidence intervals using t-test distribution
100 # degrees of freedom = n-1 = 24 - should be >30
101 t.val <- qt(alphaVal/2, df = n-1, lower.tail = FALSE)
102
103 CI_lower <- iqMean - t.val * iqse
104 CI_upper <- iqMean + t.val * iqse
105
106 # calculate using t-test
107 cat(str_glue("{(1-alphaVal)*100}% Confidence Intervals, two-sided t-test"))
108 matrix(c(CI_lower, CI_upper), ncol = 2,
109        dimnames = list("",c("Lower", "Upper")))
110
111 # t-test results in (slightly) wider confidence interval
112
113 # Check our working
114 #t.test(iqData, conf.level = 1-alphaVal, alternative = "two.sided")
115
116 cat("Our Confidence interval for the IQ of the students in the sample is: ")
117 cat(str_glue("  {round(CI_lower,2)} < mean IQ < {round(CI_upper,2)} "))
118 cat(str_glue("with a confidence level of {(1-alphaVal)*100}%"))
119
120
121 ##------------------------------------------------------------------------
122 ## Hypothesis Testing
123 # Wrangling our data
124 class(iqData) # What class of vector is our IQ variable? - numeric
125
126 # Hypothesis test:
127 # H0 : average IQ of students in school is less than or equal to  national
        average
128 # Ha : average IQ of students in school is greater than national average
129
130 # alpha = 0.05, 1-tail test, single population
131
132 # don't have variance of population, only have mean to compare against
133
134 # Test our hypothesis
135 alphaVal <- 0.05
136 popMean <- 100
137 #get test statistic
138 testStatistic <- (iqMean - popMean) / iqse
139
140 pValue <- pnorm(-abs(testStatistic))
141
142 # calculate t-test p-value
143 t_pValue <- pt(abs(testStatistic), df = n-1, lower.tail = FALSE)
144
145 matrix(c(testStatistic, n-1, t_pValue  ), ncol = 3,
146        dimnames = list("",c("t", "df", "p-value")))
147 cat(str_glue("p-value for normal distribution is {round(pValue,3)}"))
148
```

```
149
150  t.test( iqData   ,
151         mu = 100, # population mean
152         var.equal = TRUE, # The default is FALSE - don't have var for popn
153         alternative = "less", # H0: sample mean > population mean
154         conf.level = .95) #
155
156
157  # How do we interpret the output?
158  # for confidence level of 95%, we cannot reject the hypothesis (p-value >
         alpha)
159
160
161  ########################
162  # Problem 2
163  ########################
164  # function to save output to a file that you can read in later to your docs
165  output_stargazer <- function(outputFile, appendVal=TRUE, ...) {
166    output <- capture.output(stargazer(...))
167    cat(paste(output, collapse = "\n"), "\n", file=outputFile, append=appendVal)
168  }
169
170  # read in expenditure data
171  expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI_
         Fall2022/main/datasets/expenditure.txt", header=T)
172  #expenditure <- read.table("../../datasets/expenditure.txt", header=T)
173
174  # State 50 states in US
175  #Y per capita expenditure on shelters/housing assistance in state
176  #X1 per capita personal income in state
177  #X2 Number of residents per 100,000 that are "financially insecure" in state
178  #X3 Number of people per thousand residing in urban areas in state
179  #Region 1=Northeast, 2= North Central, 3= South, 4=West
180  data_headers <- c("State", "$ExpenditurePC", "$IncomePC", "FInsecureResidents"
         ,
181               "UrbanResidents", "Region")
182  regions <- c("Northeast","North Central", "South", "West")
183  names(expenditure)
184  #colnames(expenditure) <- data_headers
185
186  # Inspect the data
187  head(expenditure)
188  str(expenditure)
189  summary(expenditure)
190
191  #investigate spending on Housing assistance
192
193  # Visualise
194  hist(expenditure$Y,
195       #breaks = 12,
196       main = "Histogram of spending on HA ",
```

```r
197        xlab = "$, per capita"
198 )
199
200 plot(density(expenditure$Y),
201        main = "PDF of spending on HA ",
202        xlab = "$, per capita"
203 )
204
205 pairs(~Y + X1 + X2 + X3, expenditure)
206
207 qqnorm(expenditure$Y)
208 qqline(expenditure$Y,
209          distribution = qnorm)
210
211
212 # create plots of Y and Xn
213 onefile <- TRUE
214 #pdf( file = if(onefile) "expenditure_plots.pdf" else "expenditure_plots%03d.
       pdf")
215 #pdf("plot_example.pdf" )
216
217 ggplot(expenditure) +
218    geom_point(aes( Y, X1), colour = "blue") +
219    geom_smooth(aes( Y, X1))
220
221 ggplot(expenditure) +
222    geom_point(aes( Y, X2), colour = "blue") +
223    geom_smooth(aes( Y, X2))
224
225 ggplot(expenditure) +
226    geom_point(aes( Y, X3), colour = "blue", ) +
227    geom_smooth(aes( Y, X3), colour = "red")
228
229
230 ggplot(expenditure) +
231    geom_point(aes( STATE, Y), colour = "green") +
232    geom_point(aes( STATE, X1), colour = "blue")
233
234 ggplot(expenditure) +
235    geom_point(aes( STATE, X1/Y), colour = "green")
236
237
238 ggplot(expenditure) +
239    geom_point(aes( Y, X1), colour = "blue")
240
241 ggplot(expenditure) +
242    geom_point(aes( Y, X2), colour = "green")
243
244 ggplot(expenditure) +
245    geom_point(aes( Y, X3)) +
246    geom_smooth(aes( Y, X3))
```

```
247
248  #main = "Income  per  capita  vs  spending  on  HA "
249  ggplot ( expenditure ) +
250    geom_point ( aes ( Y, X1) , colour = "blue") +
251    geom_smooth ( aes ( Y, X1 ) )
252
253  #dev. off ()  #  close  pdf  file
254
255
256  #  regional  expenditure  on  housing  assistance
257
258  ggplot ( expenditure ) +
259    geom_point ( aes ( Y, X2, colour = factor ( Region ) ) ) +
260    geom_smooth ( aes ( Y, X2 ) )
261  #logarithmic  scale
262
263  #factor ( expenditure$Region ) <- regions
264  ggplot ( expenditure ) +
265    geom_point ( aes (  Region , Y, colours = factor ( Region ) ) )
266  #  more  spread  in  r4 ,  least  in  r2
267
268  #  can  see  eg  that  no  crossover  in  interquartile  ranges
269  boxplot ( expenditure$Y ~ expenditure$Region ,  #  here  we  use  formula  notation  to
          group
270            main = "Boxplot  of  per  capita  spending  on  HA  by  Region" ,
271            names=regions ,
272            ylab = "$" ,
273            xlab = "" )
274
275
276  regional_mean_table <-expenditure %>% #  Tidyverse  method  for  grouping
277    group_by ( Region ) %>%
278    summarise ( mean = mean(Y) )
279
280  regional_mean_table <- cbind ( regional_mean_table , regions )
281
282  output_stargazer ("regional_means.tex" , appendVal = FALSE, regional_mean_table
          [ , -1])   #  file  fragment
283
284  ggplot ( re_means ) +
285    geom_point ( aes ( regions , mean, colour = regions ) , size=3)
286
287  #  West  region  has  highest  per  capita  mean  expenditure  on  housing  assistance
288
289  #————————————————————————————————————————————————————————
290  #  look  at  income  vs  expenditure  on  HA,  by  region
291  #factor ( expenditure$Region ) <- regions
292  ggplot ( expenditure ) +
293    geom_point ( aes ( Y, X1) ) +
294    geom_smooth ( aes (Y, X1 ) )
295
```

```r
296 ggplot(expenditure) +
297   geom_point(aes( Y, X1, colour= regions[Region], shape= regions[Region]))
298
299 ggplot(expenditure) +
300   geom_point(aes( Y, X1, colour= regions[Region], shape= regions[Region])) +
301   geom_smooth(aes(Y, X1))
302
303 ggplot(expenditure) +
304   geom_point(aes( Y, X1, colour= factor(Region), shape= factor(Region))) +
305   geom_smooth(aes( Y, X1, colour = factor(Region)))
306
307
308
309
310 ggplot(data = expenditure) +
311   geom_point(mapping = aes(x = Y, y = X1)) +
312   facet_wrap(~ Region, nrow = 2)
313
314
315 ##  try - todo
316 mat <- as.matrix(with(expenditure, table(Y, Region)))
317
318
319 barplot(height = mat,
320         beside = TRUE,
321         legend.text = TRUE,
322         args.legend = list(x = "topleft",
323                            cex = 0.4,
324                            box.col = "white"))
325
326
327 # run an example regression, to show how to save table
328
329
330 lm(Y~X1, data=expenditure)
331 lm(Y~X2, data=expenditure)
332 lm(Y~X3, data=expenditure)
333
334 # execute function and check ls() to make sure it worked
335 ls()
```