

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 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)
```

1. Find a 90% confidence interval for the average student IQ in the school.
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 $\alpha = 0.05$.

```

1 # Q1: Find 90% confidence interval for the average student IQ in the school
2
3 # As stated in Logan (n.d.), a confidence interval is "the interval that
  contains
4 # the population parameter with probability 1 - ." In other words, the CI is
  the
5 # range of values within which the relevant population statistic is expected
  to sit
6 # within - to the stated level of certainty. This is calculated using values
  from
7 # a randomly-selected sample.
8
9 # If calculating manually for the t-test method, this requires finding the
  sample
10 # mean, it's standard deviation and then standard error for this mean, the
11 # corresponding t-score for the chosen confidence level, and the margin of
  error.
12 # Link: Logan (n.d.) https://bookdown.org/logan\_kelly/r\_practice/p09.html
13
14 # There are two ways of finding the 90% confidence interval - firstly, through
15 # using the t-distribution and manually calculating each component value
  before
16 # inputting these into the qt() function:
17
18 # Find the sample mean ie the central tendency
19 mean_y <- mean(y)
20
21 # Find the standard error for the population statistic
22 n_y <- length(y) # find the n-size of the sample
23 standard_deviation_y <- sd(y) # also possible to calculate this manually by
  finding
24 # the sample variability using var() and finding its squareroot ie sqrt(var(y)
  )
25 standard_error_for_y <- standard_deviation_y/sqrt(n_y)
26
27 # Find the corresponding t-score for 90% confidence interval - we are not
  assuming
28 # direction, and so we are using a two-tail test
29 confidence_interval <- 0.1
30 degrees_of_freedom <- n_y - 1
31 t_score <- qt(p = confidence_interval/2, df = degrees_of_freedom, lower.tail =
  FALSE) # where
32 # p is the vector of probabilities for each tail, df is one less the N-size as
  we
33 # exclude the mean_y value, and we indicate we are using two-tail test
34
35 # Find the margin of error
36 margin_of_error_y <- t_score * standard_error_for_y
37
38 # Find the confidence interval
39 lower_bound <- mean_y - margin_of_error_y

```

```

40 upper_bound <- mean_y + margin_of_error_y
41
42 # Print the 90% confidence intervals and sample y mean, using round() to give
43 # each value to two decimal places.
44 CI_values <- c(round(lower_bound, 2), round(mean_y, 2), round(upper_bound, 2))
45 names(CI_values) <- c("lower-bound", "sample mean", "upper-bound")
46 print(CI_values)
47
48 # The same can be approximated using Z-scores instead. The corresponding Z-
   score
49 # to the 90% confidence interval is 1.64.
50
51 lower_bound_z <- mean_y - (1.64 * standard_deviation_y / sqrt(n_y))
52 upper_bound_z <- mean_y + (1.64 * standard_deviation_y / sqrt(n_y))
53 # And now printing these values
54 CI_values_z <- c(round(lower_bound_z, 2), round(mean_y, 2), round(upper_bound_
   z, 2))
55 names(CI_values_z) <- c("lower-bound", "sample mean", "upper-bound")
56 print(CI_values_z)
57
58 # Note the values are similar, though the approximated confidence interval is
59 # slightly reduced (due to the Z-score slightly under-shooting the true 90%
   threshold)

```

```

1 # Q2: Next, the school counselor was curious whether the average student IQ
   in
2 # her school is higher than the average IQ score (100) among all the schools
   in
3 # the country. Using the same sample, conduct the appropriate hypothesis test
4 # with  $\alpha = 0.05$ .
5
6 # Now we are conducting a one-tail significance test for the difference
   between two
7 # means in order to reject the null hypothesis that mean_y is equal to or
   lower than
8 # mean_population
9
10 mean_population <- 100
11 t.test(x = y, mu = mean_population, alternative = "greater")
12
13 # The result indicates that  $p > 0.05$ , failing to reject the null hypothesis ie
14 # we fail to reject that the class's average IQ score is equal to or lower
   than
15 # the average IG score among all schools in the country.

```

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 # inputting these into the qt() function:
```

- 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.