

Problem Set 3

Quantitative Political Methodology (U25 363)

Due: April 3, 2018

Instructions

- Put your name at the top of your written document. Please show your work if possible. 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 have plots, attach them as well within your written document. Make sure you label clearly which question the codes correspond to. If you are not sure if work needs to be shown for a particular problem, please ask me.
- Your homework should be submitted electronically on the course GitHub page.
- This problem set is due before the beginning of class on Wednesday April 3, 2019. No late assignments will be accepted.
- Total available points for this homework is 100.

Question 1 (5 points)

Using data on the 2008 New Hampshire Democratic Party Primary, visualize the relationship between the proportion of voters for Howard Dean in the 2004 Democratic primary and the proportion of voters for Barack Obama in the 2008 Democratic primary. To get the dataset, type:

```
install.packages("faraway")  
library("faraway")  
data("newhamp")  
help("newhamp")
```

In addition to the relationship between the support for Dean in 2004 and the support for Obama in 2008, we are also interested in whether two different voting systems — hand-counted and machine-counted ballots — matter. At a minimum, you have to do the following things in a single plot:

- Properly label titles and axes
- Set the ranges of the axes appropriately
- Use different colors and symbols to indicate the two ballot systems
- Include a legend that explains colors and symbols

Question 2 (10 points)

The shape of the t-distribution varies by a parameter call “degrees of freedom,” or df for short.

- (a) When df is large, the t-distribution approximates what other distribution?
- (b) Use R to plot the standard normal distribution as well as three t-distributions with $df = 20$, $df = 3$, and $df = 1$. Present all plots on *the same set of axes*, print, and attach to your submitted homework. Also attach the code used to produce the plot. Give your plot a meaningful title and label your axes. Use a different color, shade of gray, or line type for each line so a reader can clearly see the difference. (You may find R’s help files useful. For example, try `?plot`, `?lines`, `?dt`.)

- (c) Describe what your plot shows about the t-distribution. With reference to your plot, explain how different sample sizes might affect your estimates of population parameters.

Question 3 (20 points)

Please find the data for this question by using the following code:

```
install.packages("Zelig")  
library("Zelig")  
data("voteincome")  
?voteincome
```

Make sure to show all your work for parts (b) and (d) either with R code you attach or by hand in the space provided. If you complete (b) and (d) in R, make sure to clearly label the code pertaining to each part of the problem.

You would like to test the hypothesis whether the average age among American voters is different from 50.

- (a) State the null hypothesis and the alternative hypothesis.

(b) Calculate the standard error, the z test-statistic and the p-value for your test.

(c) What is your conclusion at $\alpha = 0.05$?

(d) Calculate the 95% confidence interval for the mean age.

(e) How are your answers for parts (c) and (d) related?

Question 4 (25 points)

Please show all work for this problem in the space provided.

A librarian would like to learn how many books her patrons purchase per year. In particular she wants to test the $H_0 : \mu = 10$ against $H_a : \mu < 10$. She randomly selects 16 patrons and asks them how many books they purchase per year. She finds that $\bar{y} = 9.5$ and $s = 1.2$.

- (a) Given the sample size and the fact that you do not know the population standard deviation, which test statistic would you use?
- (b) What additional assumption do you need to use the test-statistic indicated in part (a)?
- (c) Calculate the test-statistic, and the p-value. What is your conclusion at significance level $\alpha = 0.05$?

- (d) Assume that you know the population standard deviation $\sigma = 1.2$. Can you use a test statistic different from the one indicated in part (a)? If so, what is that test-statistic called?
- (e) What assumption (if any) do you need to use the test-statistic indicated in part (d)?
- (f) Calculate the standard error, the test-statistic and the p-value. What is your conclusion at significance level $\alpha = 0.05$?
- (g) Compare your conclusions in parts (c) and (f). Explain the difference (if any) in the conclusions.

Question 5 (5 points)

A recent poll of 698 decided voters in Pennsylvania showed 341 preferred Donald Trump and 357 preferred Hillary Clinton. Let π be the population proportion of decided Pennsylvania voters who prefer Trump.

- (a) If the voters are only given two options (Trump or Clinton) and the sample size of your survey is relatively large, what type of distribution is the population distribution? What type of distribution is the sampling distribution of your survey?

- (b) What is the value of $\hat{\pi}$, the estimate of π obtained from the survey?

- (c) What is the standard error of this estimate?

- (d) Give the 95% confidence interval for the value of π .

Question 6 (10 points)

Field experiments have become an important tool to understand various political phenomena. For example, “Getting Out the Vote in Local Elections: Results from Six Door-to-Door Canvassing Experiments” by Green, Gerbern, and Nickerson (2003) is an early field experimental work on political behavior. The paper is available at <http://onlinelibrary.wiley.com/doi/10.1111/1468-2508.t01-1-00126/full>, but it’s also in the class readings folder. Read the abstract (and introduction if you wish) of the paper and answer the following questions.

- (a) As succinctly as possible, what is the causal claim being made by the authors?

- (b) What is the “treatment” (or predictor) variable?

- (c) What is the outcome variable?

- (d) What allows the authors to claim that their findings are causal?

Question 7 (15 points)

For the 2006 GSS, a comparison of males and females on the number of hours a day that the subject watched TV gave:

Group	N	Mean	St.Dev	SE Mean
Females	1117	2.99	2.34	0.070
Males	870	2.86	2.22	0.075

- (a) Conduct all parts of a significance test to analyze whether the population means differ for females and males. Interpret the p-value, and report the conclusion for α -level = 0.05.
- (b) If you were to construct a 95% confidence interval comparing the means, would it contain 0? (You can answer based on the result in (a), without finding the interval.)

- (c) Do you think that the distribution of TV watching is approximately normal? Why or why not? Does this affect the validity of your inferences? Explain your answer.

Question 8 (10 points)

Imagine that the data above is changed as below (note the changed sample size). A comparison of males and females on the number of hours a day that the subject watched TV gave:

Group	N	Mean	St.Dev	SE Mean
Females	11	2.99	2.34	0.070
Males	16	2.86	2.22	0.075

Conduct all parts of a significance test to analyze whether the population means differ for females and males. Interpret the p-value, and report the conclusion for α -level = 0.05.