

WDS, HW 1

Group Member 1

Group Member 2

Group Member 3

Group Member 4

Group Member 5

Due: 10:00PM, July 13, 2021

Contents

1	Overview	2
1.1	Objectives	2
1.2	Instructions	2
2	Review materials	2
3	Case study: Audience Size	3
3.1	Data preparation	3
3.2	Sample properties	3
3.3	Final estimate	4

1 Overview

This is a fast-paced course that covers a lot of material. There will be a large amount of references. It is impossible to learn data science without getting your hands dirty. Please budget your time evenly. Last-minute work ethic will not work here.

Homework in this program is different from your usual homework assignment as a typical student. Most of the time, they are built over real case studies. While you will be applying methods covered in lectures, we often do not ask very specific questions. The focus will be always on the goals of the study, the usefulness of the data gathered, and the limitations in any conclusions you may draw. Always try to challenge your data analysis in a critical way. Frequently, there are no unique solutions.

1.1 Objectives

- Get familiar with R-studio and RMarkdown
- Hands-on R
- Learn data science essentials
 - gather data
 - clean data
 - summarize data
 - display data
 - conclusion
- Packages
 - dplyr
 - ggplot

1.2 Instructions

- **Homework assignments are done in a group consisting of 5 members.**
- **All work submitted should be completed in the R Markdown format.** You can find a cheat sheet for R Markdown [here](#).
- **Submit the following files, one submission for each group:** (1) Rmd file, (2) a compiled HTML (or a pdf which might require extra work) version, and (3) all necessary data files if different from our source data. You may directly edit this .rmd file to add your answers. We encourage that you at least attempt each problem by yourself before working with your teammates. Additionally, ensure that you can ‘knit’ or compile your Rmd file. It is also likely that you need to configure Rstudio to properly convert files to PDF. [These instructions](#) might be helpful.
- In general, be as concise as possible while giving a fully complete answer to each question. All necessary datasets are available in this homework folder on Canvas. Make sure to document your code with comments (written on separate lines in a code chunk using a hashtag # before the comment) so the teaching fellows can follow along. R Markdown is particularly useful because it follows a ‘stream of consciousness’ approach: as you write code in a code chunk, make sure to explain what you are doing outside of the chunk.
- A few good or solicited submissions will be used as sample solutions. When those are released, make sure to compare your answers and understand the solutions.

2 Review materials

- Study Module 1: DataPreparationEDA_WDS

- Be able to complete DataPreparationEDA_WDS.rmd

3 Case study: Audience Size

How successful is the Wharton Talk Show [Business Radio Powered by the Wharton School](#)

Background: Have you ever listened to [SiriusXM](#)? Do you know there is a **Talk Show** run by Wharton professors in Sirius Radio? Wharton launched a talk show called [Business Radio Powered by the Wharton School](#) through the Sirius Radio station in January of 2014. Within a short period of time the general reaction seemed to be overwhelmingly positive. To find out the audience size for the show, we designed a survey and collected a data set via MTURK in May of 2014. Our goal was to **estimate the audience size**. There were 51.6 million Sirius Radio listeners then. One approach is to estimate the proportion of the Wharton listeners to that of the Sirius listeners, p , so that we will come up with an audience size estimate of approximately 51.6 million times p .

To do so, we launched a survey via Amazon Mechanical Turk ([MTurk](#)) on May 24, 2014 at an offered price of \$0.10 for each answered survey. We set it to be run for 6 days with a target maximum sample size of 2000 as our goal. Most of the observations came in within the first two days. The main questions of interest are “Have you ever listened to Sirius Radio” and “Have you ever listened to Sirius Business Radio by Wharton?”. A few demographic features used as control variables were also collected; these include Gender, Age and Household Income.

We requested that only people in United States answer the questions. Each person can only fill in the questionnaire once to avoid duplicates. Aside from these restrictions, we opened the survey to everyone in MTurk with a hope that the sample would be more randomly chosen.

The raw data is stored as `Survey_results_final.csv` on Canvas.

3.1 Data preparation

- i. We need to clean and select only the variables of interest.

Select only the variables Age, Gender, Education Level, Household Income in 2013, Sirius Listener?, Wharton Listener? and Time used to finish the survey.

Change the variable names to be “age”, “gender”, “education”, “income”, “sirius”, “wharton”, “worktime”.

- ii. Handle missing/wrongly filled values of the selected variables

As in real world data with user input, the data is incomplete, with missing values, and has incorrect responses. There is no general rule for dealing with these problems beyond “use common sense.” In whatever case, explain what the problems were and how you addressed them. Be sure to explain your rationale for your chosen methods of handling issues with the data. Do not use Excel for this, however tempting it might be.

- iii. Brief summary

Write a brief report to summarize all the variables collected. Include both summary statistics (including sample size) and graphical displays such as histograms or bar charts where appropriate. Comment on what you have found from this sample. (For example - it’s very interesting to think about why would one work for a job that pays only 10cents/each survey? Who are those survey workers? The answer may be interesting even if it may not directly relate to our goal.)

3.2 Sample properties

The population from which the sample is drawn determines where the results of our analysis can be applied or generalized. We include some basic demographic information for the purpose of identifying sample bias, if

any exists. Combine our data and the general population distribution in age, gender and income to try to characterize our sample on hand.

- i. Does this sample appear to be a random sample from the general population of the USA?

Note: You can not provide evidence by simply looking at our data here. For example, you need to find distribution of education in our age group in US to see if the two groups match in distribution. Please do not spend too much time gathering evidence.

3.3 Final estimate

Give a final estimate of the Wharton audience size in January 2014. Assume that the sample is a random sample of the MTURK population, and that the proportion of Wharton listeners vs. Sirius listeners in the general population is the same as that in the MTURK population. Write a brief executive summary to summarize your findings and how you came to that conclusion.

To be specific, you should include:

1. Goal of the study
2. Method used: data gathering, estimation methods
3. Findings
4. Limitations of the study.