

Problem Set 1

Intro to R

[YOUR NAME]

Due Date: 2024-07-04

Getting Set Up

Open RStudio and create a new RMarkdown file (.Rmd) by going to File -> New File -> R Markdown.... Accept defaults and save this file as [LAST NAME]_ps1.Rmd to your code folder.

Copy and paste the contents of this .Rmd file into your [LAST NAME]_ps1.Rmd file. Then change the author: [Your Name] to your name.

We will be using the sc_debt.Rds file from the course github page (https://github.com/jbisbee1/ISP_Data_Science_2024/blob/main/data/sc_debt.Rds).

All of the following questions should be answered in this .Rmd file. There are code chunks with incomplete code that need to be filled in.

This problem set is worth 10 total points, plus two extra credit points. The point values for each question are indicated in brackets below. To receive full credit, you must have the correct code. In addition, some questions ask you to provide a written response in addition to the code.

You are free to rely on whatever resources you need to complete this problem set, including lecture notes, lecture presentations, Google, your classmates...you name it. However, the final submission must be complete by you. There are no group assignments. To submit, email the knitted output to Eun Ji Kim (kej990804@snu.ac.kr) (<mailto:kej990804@snu.ac.kr>) as a PDF by the start of class on Thursday, July 4th. If you need help converting to a PDF, see this tutorial (https://github.com/jbisbee1/ISP_Data_Science_2024/blob/main/Psets/ISP_pset_0_HELPER.pdf).

Good luck!

*Copy the link to ChatGPT you used here: _____

Question 0 [0 points]

Require tidyverse and load the sc_debt.Rds data by assigning it to an object named df.

```
require() # Load tidyverse
```

```
## Loading required package:
```

```
df <- read_rds() # Load the dataset directly from github
```

```
## Error in read_rds(): could not find function "read_rds"
```

Question 1 [1 point]

Which school has the lowest admission rate (*adm_rate*) and which state is it in (*stabbr*)?

```
df %>%  
  arrange() %>% # Arrange by the admission rate  
  select() # Select the school name, the admission rate, and the state
```

```
## Error in df %>% arrange() %>% select(): could not find function "%>%"
```

Write answer here

Question 2 [1 point]

Which are the top 10 schools by average SAT score (*sat_avg*)?

```
df %>%  
  arrange() %>% # arrange by SAT scores in descending order  
  select() %>% # Select the school name and SAT score  
  print() # Print the first 12 rows (hint: there is a tie)
```

```
## Error in df %>% arrange() %>% select() %>% print(): could not find function "%>%"
```

Write answer here

Question 3 [1 point]

Create a new variable called *adm_rate_pct* which is the admissions rate multiplied by 100 to convert from a 0-to-1 decimal to a 0-to-100 percentage point.

```
df <- df %>% # Use the object assignment operator to overwrite the df object  
  mutate() # Create the new variable adm_rate_pct
```

```
## Error in df %>% mutate(): could not find function "%>%"
```

Question 4 [1 point]

Calculate the average SAT score and median earnings of recent graduates by state.

```
df %>%  
  group_by() %>% # Calculate state-by-state with group_by()  
  summarise(sat_avg = , # Summarise the average SAT  
            earn_avg = ) # Summarise the average earnings
```

```
## Error in df %>% group_by() %>% summarise(sat_avg = , earn_avg = ): could not find function "%>%"
```

Question 5 [1 points]

Research Question: Do students who graduate from smaller schools (i.e., schools with smaller student bodies) make more money in their future careers? Before looking at the data, write out what you think the answer is, and explain why you think so.

Write a few sentences here.

Question 6 [2 points]

Based on this research question, what is the outcome / dependent / Y variable and what is the explanatory / independent / X variable? Create the scatterplot of the data based on this answer, along with a line of best fit. Is your answer to the research question supported?

```
df %>%  
  ggplot(aes(x = , # Put the explanatory variable on the x-axis  
             y = )) + # Put the outcome variable on the y-axis  
  geom_point() + # Create a scatterplot  
  geom_smooth() + # Add line of best fit  
  labs(title = '', # give the plot meaningful labels to help the viewer understand it  
        x = '',  
        y = '')
```

```
## Error in df %>% ggplot(aes(x = , y = )): could not find function "%>%"
```

Write a few sentences here.

Question 7 [2 points]

Does this relationship change by whether the school is a research university? Using the filter() function, create two versions of the plot, one for research universities and the other for non-research universities.

```
df %>%  
  filter() %>% # Filter to non-research universities  
  ggplot(aes(x = , # Put the explanatory variable on the x-axis  
             y = )) + # Put the outcome variable on the y-axis  
  geom_point() + # Create a scatterplot  
  geom_smooth() + # Add line of best fit  
  labs(title = '', # give the plot meaningful labels to help the viewer understand it  
        subtitle = '',  
        x = '',  
        y = '')
```

```
## Error in df %>% filter() %>% ggplot(aes(x = , y = )): could not find function "%>%"
```

```
df %>%  
  filter() %>% # Filter to research universities  
  ggplot(aes(x = , # Put the explanatory variable on the x-axis  
             y = )) + # Put the outcome variable on the y-axis  
  geom_point() + # Create a scatterplot  
  geom_smooth() + # Add line of best fit  
  labs(title = '', # give the plot meaningful labels to help the viewer understand it  
        subtitle = '',  
        x = '',  
        y = '')
```

```
## Error in df %>% filter() %>% ggplot(aes(x = , y = )): could not find function "%>%"
```

Question 8 [1 point]

Instead of creating two separate plots, color the points by whether the school is a research university. To do this, you first need to modify the `research_u` variable to be categorical (it is currently stored as numeric). To do this, use the `mutate` command with `ifelse()` to create a new variable called `research_u_cat` which is either “Research” if `research_u` is equal to 1, and “Non-Research” otherwise.

```
df <- df %>%  
  mutate(research_u_cat = ifelse()) # Create a labeled version of the research_u variable
```

```
## Error in df %>% mutate(research_u_cat = ifelse()): could not find function "%>%"
```

```
df %>%  
  ggplot(aes(x = , # Put the explanatory variable on the x-axis  
             y = , # Put the outcome variable on the y-axis  
             color = )) + # Color the points by the new variable you created above  
  geom_point() + # Create a scatterplot  
  geom_smooth() + # Add line of best fit  
  labs(title = '', # give the plot meaningful labels to help the viewer understand it  
        x = '',  
        color = '',  
        y = '')
```

```
## Error in df %>% ggplot(aes(x = , y = , color = )): could not find function "%>%"
```

Extra Credit [2 points]

Write a short paragraph discussing your findings. What do you think is going on in these data?

Write a few sentences here