

Problem Set 4

Regression Part 1

[YOUR NAME]

Due Date: 2024-07-19

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 `[YOUR NAME]_ps4.Rmd` to your `code` folder.

Copy and paste the contents of this `.Rmd` file into your `[YOUR NAME]_ps4.Rmd` file. Then change the `author: [Your Name]` on line 2 to your name.

We will be using a new dataset called `youtube_individual.rds` which can be found on the course github page (https://github.com/jbisbee1/ISP_Data_Science_2024/raw/main/data/youtube_individual.rds). The codebook for this dataset is produced below. All ideology measures are coded such that negative values indicate more liberal content and positive values indicate more conservative content.

Name	Description
Responseld	A unique code for each respondent to the survey
ideo_recommendation	The average ideology of all recommendations shown to the respondent
ideo_current	The average ideology of all current videos the respondent was watching when they were shown recommendations
ideo_watch	The average ideology of all videos the respondent has ever watched on YouTube (their "watch history")
nReccs	The total number of recommendations the respondent was shown during the survey
YOB	The year the respondent was born
education	The respondent's highest level of education
gender	The respondent's gender
income	The respondent's total household income
party_id	The respondent's self-reported partisanship
ideology	The respondent's self-reported ideology
race	The respondent's race
age	The respondent's age at the time of the survey

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 questions, each worth **two** points. In addition, there are additional opportunities for extra credit totaling **another 4 points**. Note that these additional EC opportunities are **very hard**. They are designed for the students who claimed that the course is easy on the midterm survey. I encourage you all to attempt all extra credit, but don't worry if you don't get the super hard ones.

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 Friday, July 12th. 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

Require `tidyverse` and load the `youtube_individual.rds` data to an object called `yt`.

```
# INSERT CODE HERE
```

Question 1 [2 points]

We are interested in how the YouTube recommendation algorithm works. These data are collected from real users, logged into their real YouTube accounts, allowing us to see who gets recommended which videos. We will investigate three research questions in this problem set:

1. What is the relationship between average ideology of recommendations shown to each user, and the average ideology of all the videos the user has watched?
2. What is the relationship between the average ideology of recommendations shown to each user, and the average ideology of the current video the user was watching when they were shown the recommendation?
3. Which of these relationships is stronger?

Start by answering all three of these research questions, and explaining your thinking.

Write answer here.

Question 2 [1 point]

Based on your previous answer, which variables are the X (predictors) and which are the Y (outcome) variables?

Write answer here.

Now create univariate visualizations of all three variables, making sure to label your plots clearly.

```
# Y: average_recommendation_ideo
yt %>%
  ggplot(aes(x = ...)) + # Put the outcome variable on the x-axis
  geom_...() + # Choose the best geom_...() to visualize based on the variable's type
  labs(x = '', # Provide clear labels to help a stranger understand!
        y = '',
        title = '')
```

```
## Error in yt %>% ggplot(aes(x = ...)): could not find function "%>%"
```

```
# X2: average_watch_ideo
yt %>%
  ggplot(aes(x = ...)) + # Put the first explanatory variable on the x-axis
  geom_...() + # Choose the best geom_...() to visualize based on the variable's type
  labs(x = '', # Provide clear labels to help a stranger understand!
        y = '',
        title = '')
```

```
## Error in yt %>% ggplot(aes(x = ...)): could not find function "%>%"
```

```
# X2: average_current_ideo
yt %>%
  ggplot(aes(x = ...)) + # Put the second explanatory variable on the x-axis
  geom_...() + # Choose the best geom_...() to visualize based on the variable's type
  labs(x = '', # Provide clear labels to help a stranger understand!
        y = '',
        title = '')
```

```
## Error in yt %>% ggplot(aes(x = ...)): could not find function "%>%"
```

Question 3 [1 point]

Let's focus on the first research question. Create a multivariate visualization of the relationship between these two variables, making sure to put the X variable on the x-axis, and the Y variable on the y-axis. Add a straight line of best fit. Does the data support your theory?

```
yt %>%
  ggplot(aes(x = ..., # Put the first explanatory variable on the x-axis
              y = ...)) + # Put the outcome variable on the y-axis
  geom_...() + # Choose the best geom_...() for visualizing this type of multivariate relationship
  geom_...() + # Add a straight line of best fit
  labs(x = '', # Provide clear labels to help a stranger understand!
        y = '',
        title = '')
```

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

Write answer here.

Question 4 [1 point]

Now run a linear regression using the `lm()` function and save the result to an object called `model_watch`.

```
model_watch <- lm(formula = ..., # Write the regression equation here (remember to use the tilde ~!)
                  data = ...) # Indicate where the data is stored here.
```

```
## Error in eval(expr, envir, enclos): '...' used in an incorrect context
```

Using either the `summary()` function (from base R) or the `tidy()` function (from the `broom` package), print the regression result.

```
# INSERT CODE HERE
```

In a few sentences, summarize the results of the regression output. This requires you to translate the statistical measures into plain English, making sure to refer to the units for both the X and Y variables. In addition, you must determine whether the regression result supports your hypothesis, and discuss your confidence in your answer, referring to the p -value.

Write answer here.

Question 5 [1 point]

Now let's do the same thing for the second research question. First, create the multivariate visualization and determine whether it is consistent with your theory.

```
yt %>%
  ggplot(aes(x = ..., # Put the second explanatory variable on the x-axis
             y = ...)) + # Put the outcome variable on the y-axis
  geom_...() + # Choose the best geom_...() for visualizing this type of multivariate relationship
  geom_...() + # Add a straight line of best fit
  labs(x = '', # Provide clear labels to help a stranger understand!
       y = '',
       title = '')
```

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

Write answer here.

Second, run a new regression and save the result to `model_current`. Then print the result using either `summary()` or `tidy()`, as before.

```
model_current <- lm(formula = ..., # Write the regression equation here (remember to use
the tilde ~!)
                    data = ...) # Indicate where the data is stored here.
```

```
## Error in eval(expr, envir, enclos): '...' used in an incorrect context
```

```
# tidy() or summary() the result to see it
```

Finally, describe the result in plain English, and interpret it in light of your hypothesis. How confident are you?

Write answer here.

EC #1 [2 points]

Based **ONLY** on the preceding analysis, are you able to answer research question 3?

Write answer here.

Question 7 [2 points + 1 EC point]

Now let's evaluate the models. Start by calculating the “mistakes” (i.e., the “errors” or the “residuals”) generated by both models and saving these as new columns (`errors_watch` and `errors_current`) in the `yt` dataset.

```
# Calculating errors
yt <- yt %>%
  mutate(preds_watch = predict(...), # Get the predicted values from the first model (Yhat)
         preds_current = predict(...)) %>% # Get the predicted values from the second model (Yhat)
  mutate(errors_watch = ..., # Calculate errors for the first model (Y - Yhat)
         errors_current = ...) # Calculate errors for the second model (Y - Yhat)
```

```
## Error in yt %>% mutate(preds_watch = predict(...), preds_current = predict(...)) %>%
: could not find function "%>%"
```

Now create two univariate visualization of these errors. Based on this result, which model looks better? Why? EC [1 point]: Plot both errors on the same graph using `pivot_longer()`.

```
# Univariate visualization of watch history model errors
yt %>%
  ggplot(aes(x = ...)) + # Put the errors from the first model on the x-axis
  geom_...() + # Choose the best geom_...() to visualize based on the variable's type
  labs(x = '', # Provide clear labels to help a stranger understand!
        y = '',
        title = '',
        subtitle = '')
```

```
## Error in yt %>% ggplot(aes(x = ...)): could not find function "%>%"
```

```
# Univariate visualization of current video model errors
yt %>%
  ggplot(aes(x = ...)) + # Put the errors from the second model on the x-axis
  geom_...() + # Choose the best geom_...() to visualize based on the variable's type
  labs(x = '', # Provide clear labels to help a stranger understand!
        y = '',
        title = '',
        subtitle = '')
```

```
## Error in yt %>% ggplot(aes(x = ...)): could not find function "%>%"
```

```
# EC [1 point]: Plot both errors on a single plot. Hint: use pivot_longer().
# INSERT CODE HERE IF YOU DARE!
```

Write answer here.

Finally, create a multivariate visualization of both sets of errors, comparing them against the X variable. Based on this result, which model looks better? Why? EC [2 points]: Create two plots side-by-side using `facet_wrap()`. This is SUPER HARD, so don't worry if you can't get it.

```
# Multivariate visualization of watch history errors
yt %>%
  ggplot(aes(x = ..., # Put the predictor on the x-axis
             y = ...)) + # Put the errors on the y-axis
  geom_...() + # Choose the best geom_...()
  geom_...() + # Add a curve line of best fit
  geom_...(...) + # Add a horizontal dashed line at zero
  labs(x = '', # Give it clear labels
        y = '',
        title = '',
        subtitle = '')
```

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

```
# Multivariate visualization of current video errors
yt %>%
  ggplot(aes(x = ...,      # Put the predictor on the x-axis
             y = ...)) + # Put the errors on the y-axis
  geom_...() + # Choose the best geom_...()
  geom_...() + # Add a curve line of best fit
  geom_...(...) + # Add a horizontal dashed line at zero
  labs(x = '', # Give it clear labels
       y = '',
       title = '',
       subtitle = '')
```

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

```
# EC [2 points]: Try to create two plots side-by-side. (SUPER HARD)
# INSERT CODE HERE IF YOU DARE!
```

Write answer here.

Question 8 [2 points]

Calculate the **Root Mean Squared Error (RMSE)** using 100-fold cross validation with a 50-50 split for both models. How bad are the first model's mistakes on average? How bad are the second model's mistakes? Which model seems better? Remember to talk about the result in terms of the range of values of the outcome variable! EC [1 point]: plot the errors by the model using `geom_boxplot()`. HINT: you'll need to use `pivot_longer()` to get the data shaped correctly.

```

set.seed(123) # Set the seed to ensure replicability
cvRes <- NULL # Instantiate an empty object to save the results
for(...) { # 100-fold cross validation
  # Create the training dataset
  train <- yt %>%
    sample_n(size = ..., # set the size equal to half of the original dataset
             replace = ...) # Make sure to NOT replace observations (unlike bootstrapping!)

  # Create the testing dataset
  test <- yt %>%
    ... (train) # Use anti_join() to make the test set contain every observation NOT in the train set

  # Estimate model 1 on the training dataset
  mTmp_watch <- lm(formula = ...,
                  data = ...)

  # Estimate model 2 on the training dataset
  mTmp_current <- lm(formula = ...,
                   data = ...)

  # Predict both models on the testing dataset
  test <- test %>%
    mutate(preds_watch = ...,
           preds_current = ...)

  # Calculate the RMSE
  answer <- test %>%
    mutate(errors_watch = ..., # calculate the errors for model 1
           errors_current = ...) %>% # calculate the errors for model 2
    mutate(se_watch = ..., # Square the errors
           se_current = ...) %>%
    summarise(mse_watch = ..., # Take the mean of the square errors
             mse_current = ...) %>%
    mutate(rmse_watch = ..., # Take the square root of the mean of the square errors
           rmse_current = ...) %>%
    mutate(cvInd = i) # Add the cross validation index

  # Save the result
  cvRes <- cvRes %>%
    bind_rows(...)
}

# Finally, calculate the RMSE value
mean(...)
mean(...)

# EC [1 point]:
# INSERT CODE HERE IF YOU DARE

```



```
## Error: <text>:3:8: unexpected ')\n## 2: cvRes <- NULL # Instantiate an empty object to save the results\n## 3: for(...)\n##      ^
```

Write answer here.

EC #2 [2 points]

Let's try including both X variables into a single model. Run the regression and evaluate the errors as described just as you did before. Then evaluate the RMSE for ALL 3 MODELS using 100-fold cross validation with an 80-20 split. Does this combined model perform better than the two separate models? Worse? Why?

```
# INSERT CODE HERE
```

Write answer here.