

# DATA 607 Final Project

Anthony Conrardy

2024-05-04

## Introduction

While the intent of this project was to link two very disparate data sources, it turned out to be much more complicated than expected. The data sources used for this project came from the student clinical evaluation Jotform data sources located as a report on their site, and from a manual extraction of the testing data from the Platinum Educational Group website. Since all of this data is individually protected by FERPA (Family Educational Rights and Privacy Act), we must be sure to protect the individual identity of any students during this project. Though the data source for the clinical rotations was imported into R through the URL below, the URL will be deactivated before presentation and publication on RPubS. The data will then be tidy, transformed, and then exported to two usable files where the student names will be replaced with anonymous identifiers. It will be those files that may be used for analysis and further investigation.

## Data Loading

The clinical rotation survey data was set up on the Jotform site through the identified link below. Each completed student clinical survey is given a unique Submission ID, which will remain in the data set while anonymous identifiers are assigned to each unique student. The written testing data can not be easily exported from the Platinum Education Group website. Each course has to be individually access and exported in Excel format. It is those individual files that were combined into a single Excel file that was used in this section below.

```
# URL of the Excel file Clinical Rotation Survey
excel_url <- "https://www.jotform.com/excel/241118758966065"

# Define the path where you want to save the downloaded file
download_path <- "D:/Documents/R_Working_Directory/DATA 607 Final Project/excel_1.xlsx"

# Download the file
download.file(excel_url, download_path, mode = "wb")

# Read the downloaded Excel file
clinical_df <- read_excel(download_path)

# File location of final exam results Platinum Testing
excel_file <- "D:/Documents/R_Working_Directory/DATA 607 Final Project/final_exam_results.xlsx"

# read the Excel file into R
testing_df <- read_excel(excel_file)
```

## Tidying and Transforming

In this section we had to clean up the files so we could match based on first and last name. All the names were shifted to lower case to assist in matching, and the testing dataset had to have the student name separated into first and last name.

```
# Changing column names to prepare for matching
clinical_df <- clinical_df |> rename(last_name = 'Last Name')
clinical_df <- clinical_df |> rename(first_name = 'First Name')

# make name columns all lower case
clinical_df$last_name <- tolower(clinical_df$last_name)
clinical_df$first_name <- tolower(clinical_df$first_name)

# Separate Name in testing_df
testing_df <- testing_df |> separate(Student_name, into = c("last_name", "first_name"))

# make name columns all lower case
testing_df$last_name <- tolower(testing_df$last_name)
testing_df$first_name <- tolower(testing_df$first_name)
```

## Matching Data Frames

In this section we matched to two different datasets on last and first name. The clinical survey dataset that had 1372 observations, and the student testing dataset which had 143 observations, is now combined into a matched dataset of 162 observations. It should be noted that this is acceptable since some students did multiple clinical rotations, and therefore it will be that indicator that we use in the data analysis section to see if the number of clinical experiences is associated with a higher final exam score.

```
# Matching Columns
matched_df <- merge(clinical_df, testing_df, by = c('last_name', 'first_name'))
```

## Anonymous Identifier and Dataset Export

In this section we assign an anonymous identifier to the matched dataframe students and then assign that identifier to the two imported data sets. Once done, we will remove the first and last names from both datasets and only keep the variables of interest assigned to the unique identifiers. We will then export those files to GitHub where they can be accessed for the analysis section of this project. It should also be noted that the original intent of the project proposal was to analyze the student sentiment and see if there was a correlation with written exam scores. Unfortunately, the data obtained had very few commentary entries to analyze and most had 5 as a likert scale response, which did not seem to hold much value in analyzing and resulting in a change of plan.

```
# Create the anonymous identifier by grouping by last name and then assigning
# an anon_id using group_indices function
Identifier_df <- matched_df |> group_by(last_name) |> mutate(anon_id = group_indices())

# Create a Look Up table based upon last name
lookup_table <- Identifier_df |> distinct(last_name) |> mutate(anon_id=group_indices())

# Create a function that assigns the look up table values to the student name in
# the two datasets.
```

```

assign_ids <- function(df) {
  df <- left_join(df, lookup_table, by = 'last_name')
  return(df)
}
# Assign the anonymous identifiers to the students in the datasets
clinical_df <- assign_ids(clinical_df)
testing_df <- assign_ids(testing_df)

# Select the variables of interest and exclude student identification information.
clinical_export <- clinical_df |> select(anon_id, `Submission ID`, `Submission Date`, Date, `Appropriate`)

testing_export <- testing_df |> select(anon_id, final_exam, final_exam_retest)

# Exporting the Excel Files out of R

```