

Module 5 Exercise - Part B

Your Name Here

The following is an exercise where the first two (data preparation) steps are already done for you. Complete steps 3 and 4 to visualize the data and perform regression analysis.

1. Modify the nlsy97 dataset.

- Initialize the packages used in this analysis.

```
library(rio)
library(tidyverse)
library(plm)
library(ggplot2)
library(ggthemes)
library(broom)
library(lmtest)
library(stargazer)
```

- Load the final nlsy97 dataset extract (taken from the end of Module 2).

```
nlsy97 <- import("nlsy97.rds")
```

- Create the following new variables:

- *logparentincome*, equal to the log of parent income.
- A variable for the highest degree completed by the mother, equal to:
 - * College if *motheredys* is greater than or equal to 14.
 - * High school if *motheredys* is between 12 and 13.
 - * Less than high school if *motheredys* is less than 12.

```
nlsy97 <- nlsy97 %>% mutate(logparentincome = log(parentincome))

nlsy97 <- nlsy97 %>% mutate(mother_degree = case_when(
  motheredys >= 14 ~ "college",
  motheredys %in% 12:13 ~ "high school",
  motheredys %in% 0:11 ~ "less than high school"
))
```

- Turn the mother's degree variable into an ordered factor.

```
nlsy97 <- nlsy97 %>%
  mutate(mother_degree = factor(mother_degree,
    levels = c("less than high school",
               "high school",
               "college"),
    ordered=TRUE))
```

2. Create an extract, nlsy97_sample, modified from nlsy97, which:

- Drops missing values in the following variables:

- *parentincome*
- *motheredys*

- *gpa*
- *highestgrade*

```
nlsy97_sample <- nlsy97 %>% drop_na(parentincome,motheredyrs,gpa,highestgrade)
# Another approach would be:
# nlsy97 <- nlsy97 %>% filter(!is.na(parentincome),
#                               !is.na(motheredyrs),
#                               !is.na(gpa),
#                               !is.na(highestgrade))
```

- Subsets the dataset for only observations where the student is 18 years old, with a GPA greater than or equal to 1.3 (D+ average), with parent income between \$5,000 and \$100,000.

```
nlsy97_sample <- nlsy97_sample %>% filter(age == 18,
                                           gpa > 1.3,
                                           parentincome > 5000,
                                           parentincome < 100000)
```

3. Create a scatterplot

Create a scatterplot with:

- The dataset is 75 randomly sampled observations from **nlsy__sample**
- Log Parental income is on the x-axis, shown from values of 8 to 12.
- GPA is on the y-axis, shown from values 1 to 4.
- The color of points based on the mother's years of education.
- The color scale set using the following:

```
scale_colour_brewer(palette = "Set1")
```

- Add a title and axis labels.

4. Perform a regression analysis of the effect of parent income on GPA

- Using the **nlsy97** dataset, first remove observations with parental income less than \$5,000 and GPA less than or equal to 1.
- First run an OLS regression, with:
 - The log of GPA as the dependent variable
 - School type and log of parent income as the independent variables
- Display the results with tidy
- Test for autocorrelation and report the statistical decision.
- Re-run the regression as fixed effects regression, including both unit and time fixed effects.
- Perform autocorrelation-robust inference using the fixed effects regression
 - Test the coefficients of the model using the Stata-style HC_1 estimation of Newey-West heteroskedasticity and autocorrelated (HAC) robust standard errors.
 - Use the function `vcovNW()` for specifying the variance method inside of `coefest()`.
 - View the results with **tidy()**