

R-work for Online Assessment

Instructions

You should work through the code below and complete it. Keep the completed code and all the resulting output. Next you should answer the questions in the online quiz. Every student will see a slightly different collection of questions (as we will randomly draw 10 questions from a pool of about 20 questions).

The questions are of four types.

- 1) Questions that merely ask you to report output from your analysis.
- 2) Some questions will ask you about R code. For example, you will see a lot of gaps (XXXX) in the code and questions may ask you how to complete the code to make the code work. Sometimes the XXXX will represent one word and on other occasions it will represent a full line (or two) of code. Other questions may ask you about the output to be produced by a particular bit of code. If you want to practice these sorts of questions you could practice on Datacamp.
- 3) The third type of questions will test your understanding of econometric issues. For example: “What is the meaning of an estimated coefficient?” “Is a particular coefficient statistically significant?”
- 4) The fourth type of question, if asked, will be on general programming issues. For example: what is the meaning of a particular error message, or, how would you search for a particular piece of information.

Preparing your workfile

We add the basic libraries needed for this week’s work:

```
library(tidyverse)    # for almost all data handling tasks
library(ggplot2)      # to produce nice graphs
library(stargazer)     # to produce nice results tables
library(haven)        # to import stata file
library(AER)          # access to HS robust standard errors
source("stargazer_HC.r") # includes the robust regression display
```

Introduction

The data are an extract from the Understanding Society Survey (formerly the British Household Survey Panel).

Data Upload - and understanding data structure

Upload the data, which are saved in a STATA datafile (extension `.dta`). There is a function which loads STATA file. It is called `read_dta` and is supplied by the `haven` package.

```
data_USoc <- XXXX("2022_USoc_extract.dta")
data_USoc <- as.data.frame(XXXX) # ensure data frame structure
names(XXXX)
```

```
## [1] "pidp" "age" "jbhrs" "paygu" "wave" "cpi" "year"
## [8] "region" "urate" "male" "race" "educ" "degree" "mfsize9"
```

Let us ensure that categorical variables are stored as **factor** variables. It is easiest to work with these in R.

```
data_USoc$region <- XXXX
data_USoc$male <- XXXX
data_USoc$degree <- XXXX
data_USoc$race <- XXXX
```

Check for which regions we have data. Either use the **levels** or the **unique** function.

```
## [1] "north east" "north west"
## [3] "yorkshire and the humber" "east midlands"
## [5] "west midlands" "east of england"
## [7] "london" "south east"
## [9] "south west" "wales"
## [11] "scotland" "northern ireland"

## [1] south east east midlands north east
## [4] north west scotland wales
## [7] east of england northern ireland london
## [10] west midlands yorkshire and the humber south west
## 12 Levels: north east north west yorkshire and the humber ... northern ireland
```

Now check which **race** and **degree** categories exist in the data.

```
## [1] "no degree" "first degree" "higher degree"
## [1] "white" "mixed" "asian" "black" "other"
```

As we defined the **male** variable as a factor it has levels **male** and **female** (check **levels(data_USoc\$male)** to confirm). It would be better to relabel the variable to **gender**.

```
names(data_USoc)[names(data_USoc) == "male"] <- "gender"
```

The pay information (**paygu**) is provided as a measure of the (usual) gross pay per month. As workers work for varying numbers of hours per week (**jbhrs**) we divide the monthly pay by the approximate monthly hours ($4 \times \text{jbhrs}$). We shall also adjust for increasing price levels (as measured by **cpi**). These two adjustments leave us with an inflation adjusted hourly wage. We call this variable **hrpay** and also calculate the natural log of this variable (**lnhrpay**).

```
data_USoc <- data_USoc %>%
  XXXX(hrpay = XXXX/(cpi/100)) %>%
  XXXX(lnhrpay = XXXX)
```

As we wanted to save these additional variables we assign the result of the operation to **data_USoc**.

We also want to use a measure of annual pay ($\text{paygu} \times 12 / (\text{cpi} / 100)$) and add this variable (**annualpay**) to the dataframe (**data_USoc**). Also add the log of this variable as a variable to the dataframe and call it **lnannualpay**.

```
data_USoc <- XXXX %>%
  XXXX
```

Let's first summarise all numerical variables in our dataset, using the **stargazer** function.

XXXX

You should find, for instance, that the mean value of the unemployment rate (`urate`) is 7.955 and the standard deviation for the `age` variable is 18.295.

For later purposes we will also need variables $age^2/100$ and $\log(age)$. We now need to create these variables (`agesq` and `lnage`) and add them to the `data_USoc` dataframe.

XXXX

```
mean(data_USoc$lnage)
```

You should find the mean of `lnage` to be 3.744307.

Another variable needed later is a variable which indicates whether a respondent has a degree. We call this variable `grad`. It should be a factor variable with two levels, `degree` and `no degree`.

```
data_USoc <- data_USoc %>%  
  mutate(grad = ifelse(degree %in% c("first degree", "higher degree"), "degree",  
                        ifelse(degree == "no degree", "no degree", NA)))  
data_USoc$grad <- as_factor(data_USoc$grad)
```

Google to understand what the `ifelse()` function does.

Let's find out how many observations we have for some of our categorical variables.

```
table(data_USoc$degree)
```

```
##  
##      no degree  first degree higher degree  
##      102644      17509      12888
```

We can also create a frequency table for two variables, here `degree` as the row variable and `gender` as the column variable.

```
table(data_USoc$degree, data_USoc$gender)
```

```
##  
##              female  male  
##  no degree      56475 46169  
##  first degree    9464  8045  
##  higher degree   6037  6851
```

Create a frequency table with `region` as the row variable and `race` as the column variable.

```
XXXX(data_USoc$XXXX, XXXX)
```

You want to see this information in proportions (by region). Run the following code:

```
options(digits = 2)  
prop.table(table(region, race), margin = 1)
```

You should receive an error message. " object 'region' not found". Fix the code!

Data cleaning

We now remove (or “drop”) unusable (or “missing”) observations from our `data_USoc` dataframe. They are those observations which have missing (NA) data for `lnhrpay` (because the individual is not working) and we will remove observations for males who are 66 years or older and females who are 61 years or older.

```
data_USoc <- data_USoc %>%
  filter(!XXXX(lnhrpay)) %>%
  filter((gender == XXXX & age < 66) | (gender == "female" & XXXX < XXXX))
```

You should end up with 56778 observations.

Estimate regression models - Version 1

We shall estimate the following regression models (mod1)

$$\lnhrpay = \beta_0 + \beta_1 \text{ age} + \beta_2 \text{ agesq} + u$$

and (mod2)

$$\lnhrpay = \alpha_0 + \alpha_1 \lnage + u$$

```
mod1 <- lm(XXXX ~ XXXX+XXXX, data = data_USoc)
mod2 <- lm(XXXX)
stargazer_HC(mod1,mod2)
```

If you have done this correctly, you will find that that your estimated constant for mod1 is 0.485.

We suggest that there are two ways of modelling the relationship between **age** and **lnhrpay** in a so-called parametric way, either as a quadratic relationship or as a logarithmic one.

As we have lots of data, there is a third more flexible approach. We do this by generating a dummy variable for every integer age (ages are reported in full years only). To do this we will first have to create an age variable which treats age as a categorical, or in R terms, a factor variable. We shall call this **age_f**.

```
data_USoc <- data_USoc %>% mutate(age_f = as.factor(age))
```

With **age_f** being a factor variable, it is now straightforward to include this factor variable into a regression. We can either include a constant (**lnhrpay ~ age_f**) which will then use **age = 16** as a base category, or we can estimate the model without a constant (**lnhrpay ~ age_f - 1**) in which case all age categories enter separately.

```
mod3 <- lm(lnhrpay ~ age_f, data = data_USoc)
mod4 <- lm(lnhrpay ~ age_f -1, data = data_USoc)
stargazer_HC(mod3,mod4)
```

We now compare the fitted values for mod1, mod2 and mod4. First we add the predicted values to the dataframe. There are several ways to achieve this and I recommend you ask Dr. Google. (Think carefully about the search terms.)

```
data_USoc$pred_mod1 <- XXXX
data_USoc$pred_mod2 <- XXXX
data_USoc$pred_mod4 <- XXXX
```

Now we plot the predicted values for the three specifications. You should also change the Axis labels to “Predicted values” for the vertical axis, “Age (in Years)” for the horizontal axis and add a title (“Predicted Regression Model”) to your picture. If you google you should find the appropriate commands. (Again, think carefully about the search terms.)

```
ggplot(data_USoc, aes(x=age,y=pred_mod4)) +
  geom_point(color = "red") +
  geom_line(aes(y=pred_mod1),color = "blue") +
  geom_line(aes(y=pred_mod2),color = "darkorange") +
  XXXX + # add code to give your plot a title
  XXXX # add code to change the axis labels
```

The fit of `mod4` is the most flexible specification as it uses a coefficient for each year. Specifications `mod1` and `mod2` models model the relationship between `age` and `lnhrpay` with one and two parameters respectively.

Estimate regression models 2

Now we will estimate a quadratic model for `annualpay` (`annualpay ~ age + agesq`) on a subsets of data in order to compare these. When you know that you will be working with different subsets of data, the best way of doing that in R is to create a new factor variable (here `subset_ind`) which allows you to separate the data accordingly.

We will create two subgroups: 1) Males with a degree and 2) Males with no degree. You may want to check the values of the `grad` variable in order to define these correctly.

```
data_USoc$subset_ind <- "none" # default group
data_USoc$subset_ind[data_USoc$gender == "male" & data_USoc$grad == "degree"] <- "Male with degree"
data_USoc$subset_ind[XXXX] <- "Male without degree" # select all males with no degree
data_USoc$subset_ind <- as.factor(data_USoc$subset_ind)
table(data_USoc$subset_ind)
```

We will want to save the model predictions and for this purpose we pre-define a variable in which we will save the predictions.

```
data_USoc$pred_mod5 <- 0 # set the prediction to 0 by default
```

Now we estimate the model for the “Male with degree” subgroup. Note that the `lm` function accepts a `subset` argument which allows you to select a subset of observations, such as the group of all males with first degree.

```
mod5_md <- lm(XXXX ~ XXXX + XXXX, data = XXXX, subset = (subset_ind == XXXX))
stargazer_HC(XXXX)
data_USoc$pred_mod5[data_USoc$subset_ind==XXXX] <- mod5_md$fitted.values
```

Now we repeat the same just for the subgroup “Males without degree”

```
mod5_mnd <- XXXX
stargazer_HC(XXXX)
data_USoc$pred_mod5[XXXX] <- mod5_mnd$fitted.values
```

Now we plot the predicted values for the two specifications.

```
ggplot(data_USoc, aes(x=age,y=pred_mod5,color = subset_ind)) +
  geom_line() +
  ggtitle("Predicted Regression Model - Model 5") +
  ylab("Predicted values") +
  xlab("Age")
```

You will see that we have the “none” category plotted as well (of course we didn’t estimate this). You could remove these data before plotting

```
# remove observations with subset_ind == "none"
data_temp <- data_USoc %>% filter(subset_ind != "none")
```

```
ggplot(data_temp, aes(x=age,y=pred_mod5,color = subset_ind)) +  
  geom_line() +  
  ggtitle("Predicted Regression Model - Model 5") +  
  ylab("Predicted values") +  
  xlab("Age")
```

END OF INSTRUCTIONS