R-work for Online Assessment, 2022/23

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.
- 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

Set the working directory

```
setwd("YOUR/WORKING/DIRECTORY")
```

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

```
library(tidyverse)  # for almost all data handling tasks
library(ggplot2)  # to produce nice graphics
library(stargazer)  # to produce nice results tables
library(haven)  # to import stata file
library(AER)  # access to HS robust standard errors
library(countrycode)

source("stargazer_HC.r")  # includes the robust regression display
```

Introduction

The data are an extract from the OECD.

Task 1: Data Upload - and understanding data structure

Upload the data, which are saved in the two csv files OECD_RoadAccidents.csv (linked from here) and OECD_Passenger_Transport.csv (linked from here).

```
data_acc <- XXXX("OECD_RoadAccidents.csv")</pre>
data_pt <- XXXX("OECD_Passenger_Transport.csv")</pre>
names(data acc)
names(data_pt)
## [1] "LOCATION"
                      "INDICATOR"
                                    "SUBJECT"
                                                  "MEASURE"
                                                                 "FREQUENCY"
## [6] "TIME"
                      "Value"
                                    "Flag Codes"
## [1] "LOCATION"
                                    "SUBJECT"
                                                                 "FREQUENCY"
                      "INDICATOR"
                                                  "MEASURE"
## [6] "TIME"
                      "Value"
                                    "Flag Codes"
```

Look at the spreadsheets and understand the data structure.

Task 2: Cleaning Accident data

Run the following lines and use the information to understand the data structure.

```
data_acc[1,]
## # A tibble: 1 x 8
```

```
## # A tibble: 1 x 8
## LOCATION INDICATOR SUBJECT MEASURE FREQUENCY TIME Value `Flag Codes`
## <chr> <chr> <chr> <chr> <chr> A tibble: 1 x 8
## 1 SVN ROADACCID DEATH NBR A 1970 620 <NA>
unique(data_acc$MEASURE)
```

```
## [1] "NBR" "1000000HAB" "1000000VEH"
unique(data_acc$SUBJECT)
```

```
## [1] "DEATH" "INJURE" "ACCIDENTCASUAL"
```

In order to understand what these different measures and subjects represent you should consult the website linked above from which the data were downloaded.

Now find all the data which relate to France (country code: FRA) and Spain (country code: ESP) in 2017.

```
## # A tibble: 10 x 8
      LOCATION INDICATOR SUBJECT
##
                                         MEASURE
                                                     FREQUENCY
                                                                TIME
                                                                         Value Flag ~1
##
      <chr>
               <chr>
                          <chr>>
                                         <chr>
                                                     <chr>>
                                                               <dbl>
                                                                         <dbl> <chr>
               ROADACCID INJURE
##
   1 FRA
                                         NBR
                                                                2017 7.34e+4 <NA>
                                                     Α
   2 FRA
               ROADACCID DEATH
                                         1000000HAB A
                                                                2017
                                                                      5.15e+0 <NA>
    3 FRA
               ROADACCID ACCIDENTCASUAL NBR
                                                                2017
                                                                      5.86e+4 <NA>
##
                                                     Α
##
    4 ESP
               ROADACCID DEATH
                                         NBR
                                                     Α
                                                                2017
                                                                      1.83e+3 <NA>
##
  5 FRA
               ROADACCID DEATH
                                         NBR
                                                     Α
                                                                2017
                                                                      3.45e+3 < NA >
##
   6 ESP
               ROADACCID DEATH
                                         1000000HAB A
                                                                      3.93e+0 < NA>
                                                                2017
##
   7 ESP
               ROADACCID ACCIDENTCASUAL NBR
                                                                2017
                                                                      1.02e+5 <NA>
##
  8 ESP
               ROADACCID INJURE
                                                                      1.39e+5 <NA>
                                         NBR
                                                     Α
                                                                2017
## 9 FRA
               ROADACCID DEATH
                                         100000VEH A
                                                                2017
                                                                      7.23e-1 <NA>
## 10 ESP
               ROADACCID DEATH
                                         100000VEH A
                                                                2017 5.38e-1 <NA>
## # ... with abbreviated variable name 1: `Flag Codes`
```

You should see 10 rows of data.

We will only keep death data and only those data which measure the number of deaths by 1000000 inhabitants.

The only data left in the Value column is the deaths per 1,000,000 inhabitants variable. We therefore can now change the name of that variable to Deaths_p1M.

```
names(XXXX)[names(data_acc)=="XXXX"] <- "Deaths_p1M"</pre>
```

The datafile contains a few variables which are now redundant or we don't need any longer. We will only keep the variables we need, LOCATION, TIME and Deaths_p1M. The code below is faulty and you need to fix it

```
data_acc <- data_acc %>% Select(LOCAION, Time, Deaths_p1M)
```

The LOCATION variable represents the three digit country code of the respective observation. Let's add proper country names to our variables. For this we use the function countrycode from the countrycode package. The available information in the datafiles is in the LOCATION variable. It has a numeric ISO-3 format (origin = "iso3c"). We want to create a new variable in both of our datasets which is called country and contains the respective full English country name. You will have to consult the help for function countrycode to figure out what the destination format should be.

Task 3: Cleaning Personal Transport data

Run the following lines and use the information to understand the data structure.

```
data_pt[1,]
## # A tibble: 1 x 8
     LOCATION INDICATOR
                          SUBJECT MEASURE FREQUENCY TIME
                                                            Value `Flag Codes`
##
     <chr>>
              <chr>
                          <chr>
                                   <chr>
                                           <chr>
                                                      <dbl>
                                                             <dbl> <chr>
## 1 AUS
              PASSTRANSP RAIL
                                   MLN_PKM A
                                                       1970 12473. <NA>
unique(data pt$MEASURE)
```

```
## [1] "MLN_PKM"
unique(data_pt$SUBJECT)
```

```
## [1] "RAIL" "ROAD" "INLAND"
```

What does MLN_PKM stand for? What do the different values in SUBJECT represent? You may need to refer to the website (see link above) to check. Note that the miles traveled are given as total distance. The information is not standardised by population size. We will do this later.

You may want to check your understanding by looking at the data for France and Spain in 2017.

```
## # A tibble: 6 x 8
##
     LOCATION INDICATOR
                          SUBJECT MEASURE FREQUENCY
                                                      TIME
                                                            Value `Flag Codes`
##
     <chr>>
              <chr>>
                          <chr>>
                                  <chr>
                                           <chr>
                                                      <dbl>
                                                             <dbl> <chr>
              PASSTRANSP RAIL
                                  MLN PKM A
## 1 FRA
                                                      2017 110568 <NA>
## 2 FRA
                                  MLN PKM A
              PASSTRANSP ROAD
                                                      2017 873037 <NA>
## 3 ESP
                                  MLN_PKM A
              PASSTRANSP RAIL
                                                            27487 <NA>
## 4 ESP
              PASSTRANSP ROAD
                                  MLN_PKM A
                                                      2017 363368 <NA>
## 5 FRA
              PASSTRANSP INLAND
                                  MLN_PKM A
                                                      2017 983605 <NA>
## 6 ESP
              PASSTRANSP INLAND
                                  MLN PKM A
                                                      2017 390855 <NA>
```

We basically have three variables here, RAIL travel kilometers (or better miles), ROAD travel and total travel (INLAND). We want to keep these three variables and display them each as a column variable. To do so we use the pivot_wider command which is part of the tidyverse package.

This is a useful function and you may want to remember that this functionality exists for your later work.

Display again the data for France and Spain in 2017 to see the difference of the datafile's setup after this operation. You should see the following:

```
## # A tibble: 2 x 9
     LOCATION INDICATOR MEASURE FREQUENCY
##
                                              TIME `Flag Codes`
                                                                           ROAD INLAND
                                                                   RAIL
##
     <chr>>
              <chr>>
                          <chr>
                                             <dbl> <chr>
                                                                   <dbl>
                                                                          <dbl> <dbl>
                                   <chr>
              PASSTRANSP MLN PKM A
                                              2017 <NA>
## 1 FRA
                                                                 110568 873037 983605
## 2 ESP
              PASSTRANSP MLN_PKM A
                                              2017 <NA>
                                                                  27487 363368 390855
```

The datafile contains a few variables which are now redundant or we don't need any longer. We will only keep the variables we need, LOCATION, TIME, RAIL, ROAD and INLAND.

```
data_pt <- XXXX %>% XXXX(LOCATION, TIME, RAIL, ROAD, INLAND)
```

Now add the country name as above for data_acc.

After doing all this you should find that your data have the following dimensions:

data_acc: 1513 obs and 4 variablesdata_pt: 2886 obs and 6 variables

Task 4: Merging dataset

We will want to merge the two datasets data_acc and data_pt together. Before merging it is useful to review the variable names of both files and the number of observations in each.

```
nrow(data_acc)
## [1] 1513
names(data_acc)
## [1] "LOCATION" "TIME" "Deaths_p1M" "country"
nrow(data_pt)
## [1] 2886
names(data_pt)
## [1] "LOCATION" "TIME" "RAIL" "ROAD" "INLAND" "country"
```

The files have three variables in common, LOCATION, TIME and country. This is the information on the basis of which we want to match the data. So for instance we will want one row of data for Germany in 2017 and that row of data should contain Deats_p1M from the data_acc dataset and RAIL, ROAD and INLAND from the data_pt dataset.

We use the merge function to achieve this. As you can see from the above information, the two datasets contain different numbers of rows. This is because not all the sources have complete information. To demonstrate that, let's look at Ukraine (UKR) and find the information for Ukraine from both datasets.

```
data_acc %>% filter(LOCATION == "UKR")
```

```
## # A tibble: 24 x 4
               TIME Deaths_p1M country
##
      LOCATION
##
                <dbl>
                            <dbl> <chr>
##
    1 UKR
                 1994
                             14.6 Ukraine
##
    2 UKR
                 1995
                             14.6 Ukraine
##
    3 UKR
                 1996
                             13.0 Ukraine
##
    4 UKR
                 1997
                             11.8 Ukraine
##
    5 UKR
                 1998
                             11.0 Ukraine
##
    6 UKR
                 1999
                             10.6 Ukraine
##
    7 UKR
                 2000
                             10.5 Ukraine
    8 UKR
                 2001
                             12.3 Ukraine
    9 UKR
                             12.4 Ukraine
##
                 2002
## 10 UKR
                 2003
                             15.0 Ukraine
## # ... with 14 more rows
data_pt %>% filter(LOCATION == "UKR")
##
  # A tibble: 31 x 6
##
      LOCATION
                 TIME
                       RAIL
                              ROAD INLAND country
##
                <dbl> <dbl> <dbl>
                                    <dbl> <chr>
      <chr>
##
    1 UKR
                 1990 76038
                                NA
                                       NA Ukraine
    2 UKR
##
                 1991 70968
                                NA
                                       NA Ukraine
##
    3 UKR
                 1992 76196
                                NA
                                       NA Ukraine
##
                 1993 75896
    4 UKR
                                NA
                                       NA Ukraine
##
    5 UKR
                 1994 70882
                                       NA Ukraine
                                NΑ
```

You should see that the information for Ukraine starts in 1990 in the data_pt datafile but only in 1994 in the data_acc datafile. You can also see that the Ukraine only has information on rail travel, but not road travel.

NA Ukraine

NA Ukraine

NA Ukraine

NA Ukraine

NA Ukraine

You could decide, as we are merging data, to only keep information which is available from both datafiles, or to keep all information (Year-country combinations) which are available in at least one of the datasets. Here we decide that we want to do the latter, i.e. keep all information for now, even if it cannot be matched from the other file. So in the above case that means that we do want the information from 1990 to 1993 from Ukraine included. In effect that means that our new dataset, data_merge should have at least 2886 rows.

Which of the following commands does that? (Note that, as the information on the basis of which we match, country and year, is saved in identically named columns in both datafiles, we do not have to use the by.x and by.y options in the merge function).

```
data_merge <- merge(data_pt,data_acc, all.x = TRUE, all.y = TRUE)
data_merge <- merge(data_pt,data_acc, all.x = TRUE, all.y = FALSE)
data_merge <- merge(data_pt,data_acc, all.x = FALSE, all.y = FALSE)
data_merge <- merge(data_pt,data_acc, all.x = FALSE, all.y = TRUE)</pre>
```

If you have done this correctly, the new datafile should have 2937 observations and 7 variables.

Task 5: Calculate country averages

1995 63752

1996 59080

1997 54433

1998 49938

1999 47600

NA

NA

NA

NΑ

NΑ

##

##

##

##

6 UKR

7 UKR

8 UKR

9 UKR

... with 21 more rows

10 UKR

We now have multiple years of data for the countries in our dataset, data_avg. The next step is to calculate the averages for our four substantial variables (Deaths_p1M, RAIL, ROAD and INLAND).

You got it right, if your result, data_avg has 60 rows and you can replicate the following information:

head(data_avg, 10)

```
## # A tibble: 10 x 6
## # Groups:
                country [10]
##
      country
                             LOCATION avg_rail avg_road avg_inland avg_deaths_p1M
      <chr>
##
                             <chr>
                                           <dbl>
                                                    <dbl>
                                                                 <dbl>
                                                                                 <dbl>
##
   1 Albania
                             ALB
                                          242.
                                                    6075.
                                                                6177.
                                                                                  9.86
                                         7765.
                                                   34058.
                                                               49191.
                                                                                 12.7
##
    2 Argentina
                             ARG
##
    3 Armenia
                             ARM
                                            44.8
                                                    1945.
                                                                 1987.
                                                                                 10.1
##
  4 Australia
                             AUS
                                        11825.
                                                  224331.
                                                              236156.
                                                                                  6.68
##
  5 Austria
                             AUT
                                         8785.
                                                   55652.
                                                               62935.
                                                                                  8.81
##
    6 Azerbaijan
                             AZE
                                         1281.
                                                   12963.
                                                               13920.
                                                                                  9.72
##
  7 Belarus
                             BLR
                                        11235.
                                                     NaN
                                                                  NaN
                                                                                 13.4
##
  8 Belgium
                             BEL
                                         8161.
                                                   98049.
                                                               106104.
                                                                                  9.51
                                                                                  8.94
## 9 Bosnia & Herzegovina BIH
                                          691.
                                                     {\tt NaN}
                                                                 NaN
## 10 Bulgaria
                             BGR
                                          4680.
                                                   24248.
                                                               31212.
                                                                                 11.2
```

As you can see, not all countries will have information on all variables.

As mentioned above, the information coming from the passenger transport file are not standardised by population size (avg_rail, avg_road and avg_inland), while avg_deaths_p1M already is. We now want to standardise the three variables coming from the passenger transport file by population size. For that we need to import a file with population information. That information (for more than 200 countries) is in CountryInfo.csv. Import that file and merge all the available country info into data_avg only keeping the 60 rows which we have in data_avg.

Note that you should match by the three letter country code which in CountryInfo.csv is labelled as countryCode. You will therefore, now have to use the by.x and by.y options in the merge function.

Now check the names of the variables in your datafile.

```
names(data_avg)
```

```
[1] "LOCATION"
##
                           "country"
                                             "avg_rail"
                                                               "avg_road"
    [5] "avg_inland"
                           "avg_deaths_p1M"
                                             "popData2019"
                                                               "continentExp"
                                             "GDPpc"
                                                               "Obese_Pcent"
##
    [9] "Land_Area_sqkm"
                           "HealthExp"
## [13] "Over_65s"
                           "Diabetis"
```

You should have 14 variables. Also confirm that you do have 60 rows of data. All the new country data are from the year 2019.

Task 6: Standardisation

We now need to calculate the standardised personal transport variables.

You get it right if you can replicate the following.

```
## LOCATION avg_rail_pp avg_road_pp avg_inland_pp
## 1 ESP 408.7568 5487.916 5897.644
## 2 FRA 1119.4566 9651.300 10770.756
```

So the average number of miles traveled by rail per person in Spain is 409 miles per year. The average number of miles traveled on road per person in France is 9651. Convince yourself that the above calculation (avg_rail_pp = avg_rail/popData2019 * 1000000) did deliver miles per person. Recall that the travel distance coming from the passenger transport datafile was measured in Millions of miles. The population is the actual population. If we didn't multiply by 1,000,000 we would get a measure of average millions of miles traveled per person.

Lastly, we will rescale the GDPpc variable. It is currently defined in USD, but to simplify later analysis we want to express it in 1000 USD. So for instance Albania's GDPpc is USD 5224, but we want to express it as 5.224 [1000 USD].

```
data_avg <- data_avg %>% mutate(GDPpc = GDPpc/1000)
```

Look at the dataset to confirm that your operation was successful.

Task 7: Creating League Tables

Create a league table of the countries in your dataset where people travel most and least by rail and road. You want to display the 10 countries with the most and the 10 countries with the least average travel per person in both categories.

Here is the code for the table for the top 10 rail travel nations.

```
task_7a <- data_avg %>%
select(country, avg_rail_pp) %>% arrange(desc(avg_rail_pp))
head(task_7a,10)
```

```
##
          country avg_rail_pp
## 1
                     2908.5449
            Japan
## 2
      Switzerland
                     1591.2104
## 3
           Latvia
                     1348.2763
## 4
           Russia
                     1318.4358
## 5
          Belarus
                     1188.9556
## 6
          Ukraine
                     1153.9672
                     1119.4566
## 7
           France
## 8
                     1065.6674
          Hungary
## 9
          Austria
                      986.9251
                      915.7539
## 10
       Kazakhstan
```

Repeat this for the Top 10 road travel nations and equally find similar tables for the nations travelling least (but non-zero amounts) by rail and road.

Task 8: Estimate regression models - Version 1

We shall estimate the following three regression models (mod1)

```
avg\_deaths\_p1M = \gamma_0 + \gamma_1 \; GDPpc + w and (mod2) avg\_deaths\_p1M = \alpha_0 + \alpha_1 \; avg\_rail\_pp + v and (mod3) avg\_deaths\_p1M = \beta_0 + \beta_1 \; avg\_road\_pp + \beta_2 \; GDPpc + u mod1 <- lm(XXXX ~ XXXX, data = data_avg) mod2 <- lm(XXXX) mod3 <- lm(XXXX) stargazer_HC(mod1,mod2,mod3,omit.stat = "f")
```

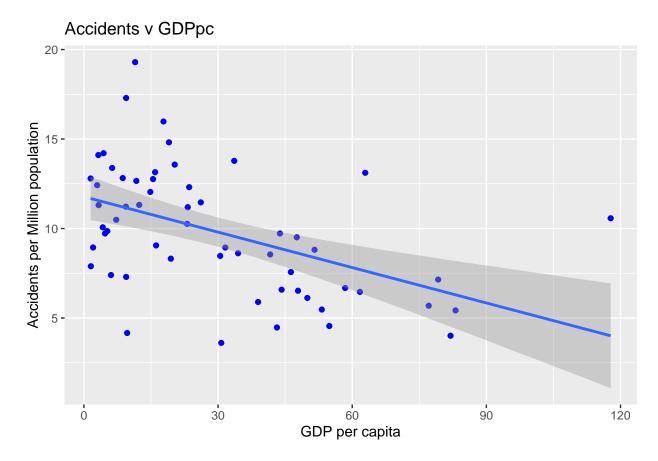
If you have done this correctly, you will find that that your estimated constant for mod1 is 11.777, the estimated constant for mod2 is 9.911 and that mod3 is estimated with 44 observations.

Task 9: Plot the data

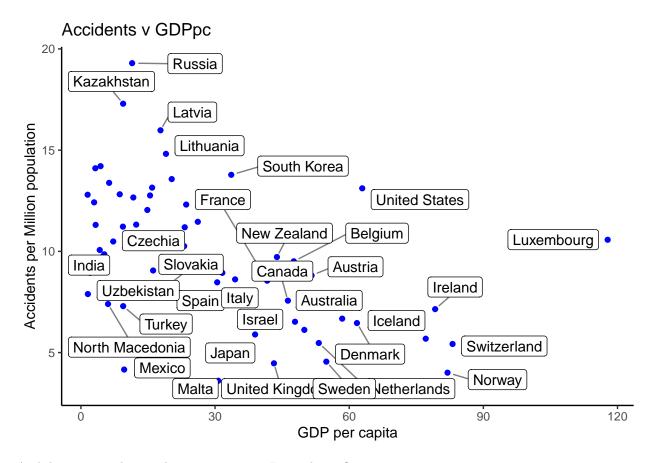
Now we want to plot the GDPpc versus the avg_deaths_p1M data in a scatter plot. Replicate the graph below.

You should change the Axis labels and add a title as shown below. If you google you should find the appropriate commands to do so. (Think carefully about the search terms.) Also, what does the geom_smooth(method='lm') line achieve? Figure that out by running the code without that line.

```
ggplot(data_avg, aes(x=XXXX,y=XXXX)) +
  geom_point(color = "blue") +
  geom_smooth(method='lm') +
  XXXX +
  XXXX +
  XXXX
```



Finally another searching challenge for you. The picture below has added data labels to the above plot. Find a way to achieve this. It is not important that the graph looks exactly like the one below, but you do want to find out how to add data labels. Dr. Google is your friend!



And does anyone know what is going on in Luxembourg?

END OF INSTRUCTIONS