

Self-Teaching Package

Version 2016

R version – Part II



Part II

Gender, employment, and wages

Overall Plan and Structure of the Exercise

The exercises in Part I demonstrated the use of household income data along with useful programming techniques for working with the LIS data. Part II emphasises the use of person-level data, including wages, demographics, and labour market information. Whereas Part I consisted entirely of calculating descriptive statistics, Part II introduces regression modelling in the final exercises.

The program that was written in the first set of exercises is now completed and can be set aside. Starting with the next exercise, you will begin the process of building up an entirely new program for Part II. Many of the techniques shown in the previous part will be useful again. In addition, users will learn how to combine LIS datasets by merging household and person files, and by concatenating multiple country-year datasets into a single file.

The general approach of the exercises is the same as in Part I. After beginning a new program in the initial exercise, each subsequent exercise will add new code to the existing program. Within each exercise, results will be produced that help to illuminate the central research themes of this section.

Research Questions

The analysis of poverty and inequality using household income, which was covered in Part 1, has always been central to research using LIS data. Over the years, however, there has been an increasing volume of work that examines individual outcomes in the labor market. The richness of the labor market data available in LIS has increased over time, and today it is possible to address many types of questions about wages and employment.

Labor market outcomes for women are one especially popular area of research. Women's rate and intensity of work shows much wider cross-country variation than men's. At the same time, on average women consistently earn lower wages.

One reflection of this growing body of research is the recent creation of the LIS *Employment Key Figures by Gender* (<http://www.lisdatacenter.org/data-access/key-figures/employment-by-gender/>), which provide information about demographics, employment, earnings, and inequality by both country and gender. In the following exercises, we will investigate several of the topics covered by the key figures, although our results will differ somewhat due to sample selection and variable coding made for the exercises.

For the exercises, we will examine three countries, using data from LIS Wave VI: the United States,

Germany, and Greece. As we will see, labor market outcomes for women show distinctly different patterns in each of these countries. Looking at persons of prime working age (which we will define as ages 25-54), our central questions will be:

- How does the percentage of prime-aged women employed in paid work vary across these three countries?
- Among those who are employed, how does the rate of *part-time* employment among women vary across the countries?
- How does employment vary by partnership and family status?
- How do wage differentials between men and women vary between across countries, across levels of educational attainment, and between immigrants and non-immigrants?

In the exercises, we will begin by producing tabulations of employment and wages for various population subgroups. In the concluding exercises, we will use linear regression to study multiple determinants of wages simultaneously, in order to better understand how family structure, education, and immigrant status are related to wages for men and for women.

Contents

1. Merging person and household files
 - Selecting a sample of prime-age workers
 - Homeownership rates by country
2. Stacking datasets
 - Employment rates by country and gender
 - Part-time employment rates by country and gender
3. Demographics and the labor market
 - Family/partnership status and employment rates
4. Employment and dependent employment
 - Non-dependent employment as a proportion of employment, by country
 - Gender wage gaps among dependent employees, by country
5. Education
 - Harmonised vs. country-specific codings
 - Gender wage gaps by country and educational attainment
6. Immigrant status
 - Understanding harmonization choices
 - Gender wage gaps by country and immigrant status
7. Multivariate analysis of wages
 - Regressing log wages on demographics, separately by gender and country
8. Pooled multi-country analysis
 - PPP adjustments for income
 - Pooled regression with multiple countries, using normalised weights

1. Merging person and household data, selecting a sample

Goal

While the exercises in Part I used only data at the household level, Part II uses data from both the household and person level files. In this exercise, we will combine the person and household files in order to create a single dataset for each country in which household data is appended to each person record.

This exercise also selects the universe of persons that we will be studying in the subsequent exercises. Since we are interested in labor market outcomes, we will restrict our attention to people of *prime age*: those who are likely to be old enough to have completed schooling and young enough to not yet be retired. In these exercises, we will define prime age persons as those between 25 and 54 years old. This is a commonly used range in statistics from the United States government and other sources, but other definitions are also possible.

Some of the variables we will be using are not always available for household members other than the head and spouse. For that reason, we will further restrict our universe to heads and partners only.

All of the variables that will be needed in the subsequent exercises are introduced here. However, for now we will only analyse one: we will summarize home ownership, which is a household-level variable, in order to measure the rate of homeownership in each of the three countries under analysis. Homeownership will be included in our later multivariate analysis, because it serves as a rough proxy for wealth, which we otherwise have no information about.

Activity

Write a program to loop through three datasets: United States 2004 (**us04**), Germany 2004 (**de04**) and Greece 2004 (**gr04**). In each country, merge the person file to the household file and keep the following variables:

- Household: Unique household identifier (**hid**) and owned/rented housing indicator (**own**)
- Person: Unique household identifier (**hid**), dataset name (**dname**), normalized person weight (**pwgt**), inflated person weight (**ppopwgt**), relationship to the household head (**relation**), partnership status (**partner**), living with own children (**children**), age (**age**), sex (**sex**), immigrant indicator (**immigr**), 3-category recoded educational attainment (**educ**), country-specific educational attainment (**educ_c**), indicator for employment (**emp**), status in employment (**status1**), indicator for part-time employment (**ptime**), and gross hourly wage in the first job (**gross1**).

Keep only those cases that are in the prime age range (between 25 and 54), and which are defined as either household heads or spouses in the variable **relation**.

Create an indicator variable equal to 1 if a person owns their house (with or without a mortgage), and 0 otherwise. Summarize this new variable to find the homeownership rate among the prime-aged persons for each country, and complete the following table:

Dataset	Homeownership %
US04	
DE04	
GR04	

Question: In which country do the largest percentage of persons of prime working age own their houses, and in which country are homeownership rates lowest?

Question: Are all the non-homeowners renters? If not, what other housing types are possible?

Guidelines

- In order to make your code easier to read, it may be helpful to store the list of variables you will be using in separate vectors, which you can refer to when you open the data file. You will need two vectors, one for household-level variables and one for person level variables. You can also store the list of datasets in a vector, and refer to it when constructing your loop.
- When loading data sets, you can use the parameter "labels" of the `read.LIS` function to tell R whether to load the data with the LIS-supplied labels, or with numeric codes. For this and subsequent exercises, we will use `labels = FALSE`.
- The simplest way to merge datasets in R is to use the `merge` command. First, load the person and household files into separate data frames, and then combine them. If you have stored the person and household files in data frames called `dp` and `dh`, you can use:

```
df <- merge(dh, dp, by = c('hid'))
```
- Remember that when recoding variables, you can find a listing of the possible values of the original variable in METIS. In this case, go to the LIS Database information. Select DE04, GR04 and US04. Select the variable `own`. Go to Cross-compare and click on the variable name to see the statistics and labels of the variables
- Like the household file, the person file contains weight variables. These variables can be used to weight by person, as an alternative to the method of multiplying household weights by number of household members that was used in the Part I. Although home ownership is a household-level variable, you will want to use the person weight to determine the proportion of *persons* who live in owner-occupied dwellings. For now, use the variable `ppopwgt`, which inflates to the total population size.

Program

```
setups <- function(ccyy) {
  # READ DATASETS
  varh <- c('hid', 'dname', 'own')
  varp <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime', 'children',
  'partner', 'status1', 'gross1', 'educ', 'immigr')
  subset <- 'age >= 25 & age <= 54 & relation <= 2200'
  dh <- read.LIS(paste(ccyy, 'h', sep = ''), labels = FALSE, vars = varh)
  dp <- read.LIS(paste(ccyy, 'p', sep = ''), labels = FALSE, vars = varp, subset = subset)
  df <- merge(dh, dp, by = c('hid'))
  # MAP NEW VARIABLES
  df$home <- ifelse(df$own %in% 100:199, 1, ifelse(df$own %in% 200:299, 0, NA))
  return(df)
}
#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
for (ccyy in datasets) {
  df <- setups(paste(ccyy, sep = ''))
  res <- round(with(df[!is.na(df$home),], sum(home*ppopwgt) / sum(ppopwgt[!is.na(home)])) *100,
  digits=2)
  print(c(ccyy, res))
}
```

Results

Dataset	Homeownership %
US04	71.5%
DE04	40.9%
GR04	67.6%

Question: In which country do the largest percentage of persons of prime working age own their houses, and in which country are homeownership rates lowest?

- Homeownership rates are highest in the United States at 71.5%, and lowest in Germany at only 40.9%.

Question: Are all the non-homeowners renters? If not, what other housing types are possible?

- No, in these three datasets there is also a category labeled “free housing”. The LIS variable definitions document explains that this can include “housing provided by employer, government or others, or illegal occupation”.

Comments

- You will notice that in this exercise the merge worked perfectly, i.e. all observations of the merging file were uniquely linked to one observation in the using file. This is always the case with LIS household and individual level files from the same dataset because all individuals belong to at least one and no more than one household.

2. Stacking data, employment rates by gender

Goal

So far, we have performed all the analysis separately for each dataset, working with only one country at a time. For this and all subsequent analyses, however, we will create a “stacked” dataset that contains information for all three countries in a single file. This means your dataset will have as many value observations as your countries altogether have stored within one file. This may offer you substantial advantages to make use of the data.

After creating a combined dataset, we will examine rates of employment and part-time employment by gender, and see how it differs among these three countries. As in the previous exercise, we will be looking only at prime-aged persons who are defined as household heads or partners of household heads.

We will be using the LIS variable **emp**, an indicator that reports whether or not a person is currently employed. This variable will ideally be coded according to the International Labor Organization (ILO) definition of employment, which identifies persons working at a certain moment in time. By this definition, a person may be considered as employed as soon as he/she has carried out any work.

Rates of employment and full-time employment among prime-age men tend to be similar and consistently high across countries, so we will be paying particular attention to differences in employment outcomes among women.

Activity

Modify your program so that it first creates a combined data file for the United States, Germany, and Greece, and then performs any necessary recoding and produces descriptive statistics.

Create a set of cross-tabulations that shows the rates of prime-age employment by gender within each country. Create another set of cross-tabulations showing the rates of part-time work by gender within each country, among those who are employed. Use your results to complete the following table below.

You should write your code so that your overall program is broken down into three subroutines. The first subroutine should contain only the code needed to create the merged, stacked dataset. The second subroutine should contain all of the data-preparation and recodings. The third subroutine should contain code that produces the summary statistics. Your overall program can then simply call these two subroutines to make the dataset and output the results. Breaking up your code in this way will be important for making the program compact and efficient in later exercises.

Dataset	Female employment rate	Part-time employment rate among employed women
US04		
DE04		
GR04		

Question: Contrast these countries in terms of their rates of female employment (high or low) and their rates of part-time employment among employed women (high or low).

Guidelines

- The `read.LIS` command can automatically stack datasets if you call it with a vector of identifiers instead of a single string:

```
df <- read.LIS(c('us04h', 'de04h', 'gr04h'))
```

This will return a single data frame containing data from the United States, Germany, and Greece. Note, however, that you cannot mix person and household file identifiers in the same call, or you will get an error. Thus you may want to first create two stacked datasets (one for person and one for household data) and then merge them. If you do this, however, note that you must merge on two variables, `dname` and `hid`, since the household identifier by itself will no longer uniquely identify cases. This is easily done by passing a vector of merging variables as the third argument of the `merge()` command.

- In base R, `weighted proportions` does not exist as such. One option is to use the function `tapply(X, INDEX, FUN, ...)` – that applies a function or operation on subset of the vector broken down by a given factor variable (or a discrete numeric variable).

For example, apply the following code to calculate the Female employment rate

```
with(df[df$sex==2 & !is.na(df$emp),], tapply(emp*ppopwgt, list(dname), sum) / tapply(ppopwgt, list(dname), sum)) * 100
```

This code generate the weighted proportion of the `emp` variable subdivided by dataset name for women. In addition, the results are multiplied by 100 to get the percentage.

Program

```
get_stack <- function(datasets, varp, varh, subset) {
  # READ DATASETS
  pp <- read.LIS(paste(datasets, 'p', sep = ''), labels = FALSE, vars = varp,
subset = subset)
  hh <- read.LIS(paste(datasets, 'h', sep = ''), labels = FALSE, vars = varh)
  df <- merge(pp, hh, by = c("dname", "hid"))
  return(df)
}
#----- RUN SCRIPTS -----
varh <- c('hid', 'dname', 'own')
varp <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime')
subset <- 'age >= 25 & age <= 54 & relation <= 2200'
datasets <- c('us04', 'de04', 'gr04')
df <- get_stack(datasets, varp, varh, subset)

'Female Employment Rate'
round(with(df[df$sex==2 & !is.na(df$emp)], tapply(emp*ppopwgt, list(dname),
sum) / tapply(ppopwgt, list(dname), sum)) *100, digits=2)

'Partime Employment rate among employed women'
round(with(df[df$sex==2 & df$emp==1 & !is.na(df$ptime)], tapply(ptime*ppopwgt,
list(dname), sum) / tapply(ppopwgt, list(dname), sum)) *100, digits=2)
```

Results

Dataset	Female employment rate	Part-time employment rate among employed women
US04	72.4%	19.7%
DE04	75.2%	45.7%
GR04	58.3%	15.8%

Question: Contrast these countries in terms of their rates of female employment (high or low) and their rates of part-time employment among employed women (high or low).

- Employment rates among prime age workers are relatively high in the United States in Germany, and lower in Greece. In the United States, most employed women work full time, while nearly half of employed German women work part time. Greece combines low employment rates with high rates of full time employment among those women who are employed.

Comments

- While **emp** will ideally contain employment status based on the ILO definition, in some cases it must be based on usual or main activity status. See the dataset-specific notes to the variable for information.
- There are several other variables that measure employment status in LIS, which may be useful for different purposes. More detailed information on current labor force status is contained in **clfs**, which is usually the variable used to construct **emp**.
- Those interested in those for whom employment is the main activity should consider the main activity variables **mainemp** and **cmas**.
- Those wishing to perform income distribution analyses according to the activity status of individuals (e.g., calculation of poverty rates for the working poor), may choose to focus on employment over a longer reference period. In this case, you need to be able to identify an individual's primary activity over the income reference period (which in LIS is normally one year), using the LIS variable **umas**.

3. Family structure and employment

Goal

In the previous exercise we examined cross-national differences in women's employment. In this exercise, we will examine the variation in employment rates among women, based on their partnership and family status. We will contrast partnered and single women. Within each of those two categories, we will contrast women without children in the household, women with young children, and women with older children. The variables created in this exercise will be useful later, when we combine family structure with other personal characteristics in a multivariate analysis of wages.

Activity

Since you already created the merged, stacked dataset in the previous exercise, you do not need to create it again. Modify your code so that the subroutine that merges and stacks the data is commented out, and add a line that simply loads the merged and stacked file at the beginning of the program.

Create a variable **achildcat**, to indicate the age of the youngest own child living in the household. This variable should be equal to 0 if there are no children under 18, equal to 1 if the youngest child is under 6 years old, and equal to 2 if the youngest child is between 6 and 17. You can create this variable based on the information in the variable **children**.

Create another variable **livepartner**, to indicate whether a person is currently living with a partner (=1) or not (=0). We will not distinguish married couples from non-married couples for this analysis, and will simply consider a person partnered in either case. If someone has a partner but is not living with them, however, we will *not* consider them partnered (**livepartner**=0).

Using the new variables you have created and the employment indicator you used in the previous exercise, produce summary statistics to fill the table below.

Employment Rates

Dataset	All women	Partnered			Single		
		No children under 18	Child under 6	Child 6-17	No children under 18	Child under 6	Child 6-17
US04							
DE04							
GR04							

Question: Which subpopulation of prime age women has the lowest employment rates in each country?

Guidelines

- You will again need to recode LIS variables, as in exercise 2.1. See that exercise for more information. In this exercises, you will be creating `achild` and `partn`.
- As in the last exercise, you can use the `tapply` function to calculate weighed proportions over subsets of the data. You can extend this command by adding additional variables to the list of categorical variables, which will allow you to analyze subsets within subsets. For example:

```
with(df[df$sex == 1 & !is.na(df$emp), ], tapply(emp * ppopwgt, list(dname, child, partn), sum) / tapply(ppopwgt , list(dname, child, partn), sum)) * 100
```

This code will produce a tabulation of employment rates, which is separated by country, by gender, and by partnership status.

Program

```
get_stack <- function(datasets, varp, varh, subset) {
  # READ DATASETS
  pp <- read.LIS(paste(datasets, 'p', sep = ''), labels = FALSE, vars = varp, subset = subset)
  hh <- read.LIS(paste(datasets, 'h', sep = ''), labels = FALSE, vars = varh)
  df <- merge(pp, hh, by = c("dname", "hid"))
  # MAP NEW VARIABLES
  df$child <- ifelse(df$children %in% c(140, 200), 0, ifelse(df$children == 110, 1,
ifelse(is.na(df$children), NA, 2)))
  df$partn <- ifelse(df$partner %in% c(100, 120, 200), 0, ifelse(is.na(df$partner), NA, 1))
  return(df)
}
#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
varh      <- c('hid', 'dname', 'own')
varp      <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime', 'children',
'partner', 'status1')
subset    <- 'age >= 25 & age <= 54 & relation <= 2200'
df        <- get_stack(datasets, varp, varh, subset)

round(with(df[df$sex == 2 & !is.na(df$emp), ], tapply(emp * ppopwgt, list(dname, child, partn),
sum) / tapply(ppopwgt, list(dname, child, partn), sum)) * 100, digits = 2)
```

Results

Employment Rates

Dataset	All women	Partnered		Single			
		No children under 18	Child under 6	Child 6-17	No children under 18	Child under 6	Child 6-17
US04	72.4%	75.9%	60.0%	73.5%	79.0%	66.6%	77.5%
DE04	75.2%	81.0%	53.4%	75.6%	86.3%	48.2%	77.7%
GR04	58.3%	50.6%	58.3%	61.4%	69.1%	64.3%	81.4%

Question: Which subpopulation of prime age women has the lowest employment rates in each country?

- In the United States, partnered women with young children have the lowest employment rates. In Greece, partnered women without children have the lowest employment rates. In Germany, however, single mothers with children have lower employment rates. This may be because of more generous child policy in Germany that makes it easier for mothers of young children to support themselves without paid employment.

Comments

- There is no clear-cut definition of a single-mother household. In this exercise, we allow other adult members to be present (as long as they are not defined as her partner). An alternative approach would be to limit the sample to households composed of a single female adult and her children. Another possibility is to limit single mother households to those with children under a specified age limit.
- When subdividing subsets of the population as has been done here, pay attention to sample sizes. In small datasets, estimates for narrowly defined groups may become very small, making estimates less reliable. The estimate for single Greek women with young children in this exercise, for example, is based on only 17 cases!

4. Dependent employment and hourly wages

Goal

In the next several exercises, we will shift from considering employment to analysing the earnings of those who are employed. We will focus our analysis on the hourly wages and thus restrict our sample to those in *dependent* employment only — that is, those who are employees. The self-employed, along with several other small categories of workers, are excluded.

In this exercise we will first determine how many workers are excluded from the analysis when the sample is restricted to those in dependent employment. We will then measure the gap in hourly wages for men and women, in each of the three countries in our study.

We will be using a measure of gross hourly wages, which is available in the three datasets we are using. In other datasets, however, it could be that only annual wages were available. In such cases, researchers must take care to account for variations in employment over the year, perhaps by restricting the sample to full-year, full-time workers.

In part I, we have introduced bottom- and topcoding as a technique to deal with extreme values. In this exercise we will also introduce a different technique that deals with extreme values: ‘trimming’ the distribution (i.e. deleting rather than recoding the extreme values). This technique is especially useful when calculating measures that are not defined for non positive values (such as logarithmic measures). In later exercises we will convert the hourly wages into logs, and thus we need to make sure that the sample that we analyse at this stage is the same that we will keep for our final analysis.

Activity

Recode the variable **status1** to create a new variable **depemp** that indicates whether a person is in dependent employment. Using this variable, produce summary statistics reporting the proportion of dependent employment among prime-age male workers and among prime-age female workers, and fill in the following table.

Employment Rates

Dataset	Men		Women	
	Dependent Employment (%)	Non-dependent employment (%)	Dependent Employment (%)	Non-dependent employment (%)
US04				
DE04				
GR04				

Next, use the LIS hourly wage variable **gross1** to construct a new hourly wage variable **hourwage**, where the bottom and the top of the distribution are corrected as follows:

- at the top, we will carry out the same topcode as used in Part I: hourly wages greater than ten times the national median are set to ten times the national median (make sure you calculate the median and apply the topcode *separately* for each country);
- at the bottom, we will 'trim' the distribution so that observations with negative and zero hourly wages are set to missing.

Using this new hourly wage variable, calculate the *gender wage gap* for dependent employees in each country. The gender wage gap is defined here as the ratio of the median wages of women to the median wages of men. Use your results to complete the table below.

	Gender wage gap for dependent employees
US04	
DE04	
GR04	

Question: Does the percentage of workers not in dependent employment differ substantially across countries? Does it differ between men and women?

Question: Which country has the most wage inequality between men and women, among dependent employees?

Guidelines

- Create a new variable for dependent employment using the recoding techniques from the earlier exercises.
- You can use the `wNtile(var, wgt, split)` function used in the part I of the self-teaching package to calculate the weighted median for the purpose of topcoding.

```
topline <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],], 10 * wNtile(hrwg,
ppopwgt, 0.5))
```

- As in the last exercise, you can use the `tapply` function to calculate the employment rates.
- Use the `wNtile(var, wgt, split)` function to calculate median wages for women and men by country with a for loop technics.

```
for (i in (1:length(datasets))) {
  for (j in 1:length(unique(df$sex))) {
    mat[j, i] <- ...
  }
}
```

- Then allocate the results to a (2LX3C) matrix and calculate the ratio of the values per line:
`round(mat[2,] / mat[1,], digits = 2)`

Program

```
get_stack <- function(datasets, varp, varh, subset) {

# READ DATASETS
pp <- read.LIS(paste(datasets,'p',sep=''),labels=FALSE, vars=varp, subset = subset)
hh <- read.LIS(paste(datasets,'h', sep=''), labels=FALSE, vars=varh)
df <- merge(pp, hh, by = c("dname", "hid"))

# MAP NEW VARIABLES
df$sex <- ifelse(df$sex == 1, 0, 1)
df$dept <- ifelse(df$status1 %in% c(100:122),1,ifelse(is.na(df$status1),NA,0))
df$hrwg <- df$gross1
df$hrwg <- ifelse(df$hrwg <= 0, NA, df$hrwg)
for (i in 1:length(datasets)) {
  topline <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],], 10 * wNtile(hrwg,
  ppopwgt, 0.5))
  df$hrwg <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],], ifelse(df$hrwg >
  topline, topline, df$hrwg))
}
return(df)
}

wNtile <- function(var, wgt, split) {
  x <- var[order(var)]
  y <- wgt[order(var)]
  z <- cumsum(y) / sum(y)
  cop <- rep(NA,length(split))
  for (i in 1:length(cop)) {
    cop[i] <- x[Find(function(h) z[h] > split[i], seq_along(z))]
  }
  return(cop)
}

#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
varh <- c('hid', 'did','dname', 'own')
varp <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime',
'children', 'partner', 'status1', 'gross1')
subset <- 'age >= 25 & age <= 54 & relation <= 2200'
df <- get_stack(datasets, varp, varh, subset)

# EMPLOYMENT RATES
print('Employment Rate')
round(with(df[df$emp == 1 & !is.na(df$dept), ], tapply(dept * ppopwgt, list(dname,
sex), sum) / tapply(ppopwgt, list(dname,sex), sum)) * 100, digits=2)

# GENDER WAGE GAP
mat <- matrix(NA, 2 ,3)
colnames(mat) <- datasets
for (i in (1:length(datasets))) {
  for (j in 1:length(unique(df$sex))) {
    mat[j,i] <- with(df[df$dname==datasets[i] & !is.na(df$hrwg) & df$sex== j-1,] ,
wNtile(hrwg, ppopwgt, 0.5))
  }
}

'Gender Wage Gap'
round(mat[2, ] / mat[1, ], digits = 2)
```

Results

Employment Rates

Dataset	Men		Women	
	Dependent Employment (%)	Non-dependent employment (%)	Dependent Employment (%)	Non-dependent employment (%)
US04	86%	14%	91%	9%
DE04	88%	12%	92%	8%
GR04	66%	34%	70%	30%

Gender wage gaps

	Gender wage gap for dependent employees
US04	0.75
DE04	0.82
GR04	0.85

Question: Does the percentage of workers not in dependent employment differ substantially across countries? Does it differ between men and women?

- Greece has a much higher rate of non-dependent employment (which is primarily self-employment). In all the countries, women have somewhat higher rates of dependent employment than men. Keep in mind, therefore, that the results in the subsequent exercises may be somewhat unrepresentative, particularly for Greece, because they exclude a substantial number of workers.

Question: Which country has the most wage inequality between men and women, among dependent employees?

- The United States shows the largest gender wage gap. Among prime age workers, the median hourly wage of women is only 75% that of men.

Comments

- The hourly wage variable **gross1** reports the wage in the first job. If information is available for a second job, it will be in the variable **gross2**. All other job characteristics are also organised in this way (e.g. industry, occupation, etc.).
- The wage gap calculated here is based on the median, but some researchers calculate an alternative version based on the mean, which will give slightly different results.

5. Hourly wages, education, and country-specific variables

Goal

This exercise continues the analysis of gender wage gaps in hourly wages among those in dependent employment, which we started to program in the previous exercise. In this exercise, we will see how gender disparities in wages differ by educational attainment.

This exercise also demonstrates the use of two different LIS variables measuring educational attainment. One is fully standardised for cross-national compatibility, but contains few categories. The other may contain more information, but has country-specific codes, and thus requires researchers to perform their own standardisation.

The standardised variable is called **educ**, which is based on the International Standard Classification of Education (ISCED). The non-standardised version (from which **educ** is constructed) is **educ_c**. This is one of many attributes for which LIS provides both a standardised and country-specific variable. Any variable ending in **_c** is non-standardised, meaning that it can have different contents in different datasets. It is important to carefully examine the dataset-specific documentation before using such variables.

Activity

Add code to your program to create a table cross-tabulating the variables **educ** and **educ_c** for each country. This will show how the standardised variable was constructed in each case. Be sure to:

- include missing values in your table, so that you can see whether any of the cases in the original education variable could not be allocated to a category in the standardised version;
- remove the value labels from the tabulation of the **educ_c** variable (since the value labels of the **_c** variables are by definition dataset-specific, in a stacked dataset with observations from several LIS datasets, the value labels of those variables will be incorrect, as they can only refer to one specific LIS dataset – usually the last one that was used to construct the stacked data, see further details in the comments section of this exercise).

Using the hourly wage variable **hourwage** that you created in the last exercise, calculate the gender wage gap by education for each country, and complete the table below. The gender wage gap is defined as it was in the previous exercise. To obtain the earnings ratio by education, simply calculate the ratio separately for individuals in each of the three categories of the standardised education variable.

Gender wage gaps by educational attainment

	Low education	Medium education	High education
US04			
DE04			
GR04			

Question: For each of the three countries, what are the categories in the original dataset that are recoded as “high education” in the standardised education variable?

Question: Are there any categories in the original **educ_c** variable that could not be translated into the standardised form?

Question: In general, which educational attainment category shows the greatest earnings inequality between genders? How do the patterns differ by country?

Guidelines

- Because the program so far has been designed to load variables without their labels, you will have to consult the codebook for the dataset at *Our Data* → *By Country (under LIS Database)* → *<country name>*, in order to determine what the country-specific education codes refer to. If you like, you can write your program to load the education variables with labels before tabulating, but this may be somewhat redundant and cumbersome because of the way the program has been designed so far.
- As in the previous exercises, you can produce tabulations of the wage ratio for each country by using **tapply()** function and loop technics.

Program

```
get_stack <- function(datasets, varp, varh, subset) {
  # READ DATASETS
  pp <- read.LIS(paste(datasets, 'p', sep = ''), labels = FALSE, vars = varp,
subset = subset)
  hh <- read.LIS(paste(datasets, 'h', sep = ''), labels = FALSE, vars = varh)
  df <- merge(pp, hh, by = c("dname", "hid"))
  # MAP NEW VARIABLES
  df$sex <- ifelse(df$sex == 1, 0, 1)
  df$hrwg <- df$gross1
  df$hrwg <- ifelse(df$hrwg <= 0, NA, df$hrwg)
  for (i in 1:length(datasets)) {
    topline <- with(df[!is.na(df$hrwg) & df$dname == datasets[i], ], 10 * wNtile(hrwg,
ppopwgt, 0.5))
    df$hrwg <- with(df[!is.na(df$hrwg) & df$dname == datasets[i], ], ifelse(df$hrwg >
topline, topline, df$hrwg))
  }
  return(df)
}
wNtile <- function(var, wgt, split) {
  x <- var[order(var)]
  y <- wgt[order(var)]
  z <- cumsum(y) / sum(y)
  cop <- rep(NA, length(split))
  for (i in 1:length(cop)) {
    cop[i] <- x[Find(function(h) z[h] > split[i], seq_along(z))]
  }
  return(cop)
}
#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
varh <- c('hid', 'dname', 'own')
varp <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime',
'children', 'partner', 'status1', 'gross1', 'educ')
subset <- 'age >= 25 & age <= 54 & relation <= 2200'
df <- get_stack(datasets, varp, varh, subset)

# GENDER WAGE GAP
ctry_list <- list()
for (k in (1:length(datasets))) {
  mat <-
matrix(NA, length(unique(df$sex[!is.na(df$sex)])), length(unique(df$educ[!is.na(df$educ)])))
  colnames(mat) <- c('Low', 'Medium', 'high')
  for (j in (1:length(unique(df$educ[!is.na(df$educ)])))) {
    for (i in (1:length(unique(df$sex[!is.na(df$sex)])))) {
      mat[i,j] <- with(df[df$dname == datasets[k] & !is.na(df$hrwg) & !is.na(df$sex) &
!is.na(df$educ) & df$sex == i - 1 & df$educ == j, ], wNtile(hrwg, ppopwgt, 0.5))
    }
  }
  ctry_list[[datasets[k]]] <- round(mat[2, ] / mat[1, ], digits = 2)
}

'Gender Wage Gap'
ctry_list
```

Results

Gender wage gaps by educational attainment

	Low education	Medium education	High education
US04	0.71	0.71	0.78
DE04	0.69	0.84	0.83
GR04	0.67	0.78	0.96

Question: For each of the three countries, what are the categories in the original dataset that are recoded as “high education” in the standardised education variable?

- In the United States, high education combines those with associate degrees, bachelor's degrees, and advanced degrees (masters, professional school, or doctorate).
- In Germany, high education combines those classified as having “higher vocational” (of any kind) or “higher education”.
- In Greece, high education includes those with tertiary graduate or postgraduate level education.

Question: Are there any categories in the original **educ_c** variable that could not be translated into the standardised form?

- In the United States, all values of **educ_c** receive a value in **educ**. In Germany, a small number of persons categorised as “inadequately, other diploma” or “still in education” are set to missing. In Greece, a small number of persons listed as “still in education” are set to missing.

Question: In general, which educational attainment category shows the greatest earnings inequality between genders? How do the patterns differ by country?

- In all three countries, there is a smaller gender wage gap among highly educated workers. This is particularly notable in Greece. In that country, wage inequality is greater among the low-educated than in the United States and Germany, but there is near equality among the highly educated.

Comments

You may have some doubts and questions why we did advise you in the exercise to not show the labels of **educ_c** for the cross-tabulation of **educ** and **educ_c**. As you are aware we did append values for several countries for each variable to get the stacked file. In your stacked file, for standardised variables these values have all the same meaning, as the values and labels are completely the standardised. However, it is more complicated for non-standardized values and labels of **_c** variables, as each dataset has its own values and own meaning, as indicated by the labels attached to the data.

Be aware that while appending the data, your programming software will very likely overwrite the label automatically. Thus we do in general advise you to drop the labels from the variables **_c** within your code. You always have the full information on the labels in the codebooks. However, if you prefer to keep the full labels somewhere in the data there are several solutions.

The simple solution is that you tabulate each country separately (see for example exercise 2 of part II) before you generate the stacked version. As a second solution, you can also easily rename the variables **_c** to the specific **_`ccyy'** of each dataset - this way you will append a separate variable **_`ccyy'** for each of your datasets, which does have only observations for the specific **`ccyy'** with the country specific labels attached. Be aware that you then need to tabulate for each **`ccyy'** separately to get the right percentage of missings!

6. Immigration and wages, understanding harmonisation

Goal

Each of the countries we are examining has a significant immigrant population, and their labor market outcomes are often very different from those of the non-immigrant population. In this exercise we will compare the wages of immigrants and non-immigrant men and women, just as we compared individuals of different educational levels in the last exercise.

LIS provides a variable indicating whether someone is an immigrant, which we will be using in this exercise. However, the choices that go into constructing this variable are complex, because the information available to construct it varies widely from country to country. It is important to understand the assumptions behind variables such as this one, because in some cases researchers may prefer to develop their own standardisation procedures based on their particular needs.

Activity

Using the top and bottom coded hourly wage variable, calculate the gender earnings ratio by immigration status for each country, and complete the table below. The gender earnings ratio is computed just as in the previous exercise, except that you will now subdivide the population into immigrants and non-immigrants, rather than by educational attainment categories.

Gender earnings ratios by immigration status

	Non-immigrants	Immigrants
US04		
DE04		
GR04		

Question: What information is used to construct the **immigr** variable? If you wanted to determine how the indicator was constructed in a particular dataset, what other variables would you need to look at?

Question: Is gender earnings inequality larger among immigrants or non-immigrants? Does this differ by country?

Guidelines

- The coding required for this exercise is essentially the same as in the previous one.

Program

```
get_stack <- function(datasets, varp, varh, subset) {
# READ DATASETS
pp <- read.LIS(paste(datasets, 'p', sep = ''), labels=FALSE, vars=varp, subset=subset)
hh <- read.LIS(paste(datasets, 'h', sep = ''), labels = FALSE, vars = varh)
df <- merge(pp, hh, by = c("dname", "hid"))

# MAP NEW VARIABLES
df$sex <- ifelse(df$sex == 1, 0, 1)
df$dept <- ifelse(df$status1 %in% c(100:122), 1, ifelse(is.na(df$status1), NA, 0))
df$hrwg <- df$gross1
df$hrwg <- ifelse(df$hrwg <= 0, NA, df$hrwg)
for (i in 1:length(datasets)) {
  topline <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],], 10*wNtile(hrwg,
  ppopwgt,0.5))
  df$hrwg <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],], ifelse(df$hrwg >
  topline, topline, df$hrwg))
}
return(df)
}
wNtile <- function(var, wgt, split) {
  x <- var[order(var)]
  y <- wgt[order(var)]
  z <- cumsum(y) / sum(y)
  cop <- rep(NA,length(split))
  for (i in 1:length(cop)) {
    cop[i] <- x[Find(function(h) z[h] > split[i], seq_along(z))]
  }
  return(cop)
}
#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
varh <- c('hid', 'dname', 'own')
varp <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime',
'children', 'partner', 'status1', 'gross1', 'educ', 'immigr')
subset <- 'age >= 25 & age <= 54 & relation <= 2200'
df <- get_stack(datasets, varp, varh, subset)

# GENDER WAGE GAP
ctry_list <- list()
for (k in (1:length(datasets))) {
  mat <- matrix(NA, length(unique(df$sex[!is.na(df$sex)])),
length(unique(df$immigr[!is.na(df$immigr)])))
colnames(mat) <- c('Non Immigrant', 'Immigrant')
for (j in (1:length(unique(df$immigr[!is.na(df$immigr)])))) {
  for (i in (1:length(unique(df$sex[!is.na(df$sex)])))) {
    mat[i,j] <- with(df[df$dname == datasets[k] & !is.na(df$hrwg) & !is.na(df$sex)
& !is.na(df$immigr) & df$sex== i-1 & df$immigr== j-1,],
wNtile(hrwg,ppopwgt,0.5))
  }
}
ctry_list[[datasets[k]]] <- round(mat[2, ] / mat[1, ], digits = 2)
}

'Gender Wage Gap'
ctry_list
```

Results

Gender earnings ratios by immigration status

	Non-immigrants	Immigrants
US04	0.76	0.82
DE04	0.82	0.74
GR04	0.89	0.82

Question: What information is used to construct the **immigr** variable? If you wanted to determine how the indicator was constructed in a particular dataset, what other variables would you need to look at?

- As explained in the “LIS Variable Definitions” document, Immigrants are defined by LIS as all persons who have that country as country of usual residence and (in order of priority):
 - whom the data provider defined as immigrants;
 - who self-define them-selves as immigrants;
 - who are the citizen/national of another country;
 - who were born in another country.
- **immigr** could be constructed out of **citizen**, **ctrybrth**, **yrsresid**, **ethnic_c** and **immigr_c**.

Question: Is gender earnings inequality larger among immigrants or non-immigrants? Does this differ by country?

- In the United States, the gender wage gap is greater among non-immigrants, but in Germany and Greece it is greater among immigrants.

Comments

- The definition of immigrant used in **immigr** may differ substantially from dataset to dataset. The variables that may be used in its construction include **citizen**, **ctrybrth**, **yrsresid**, **ethnic_c** and **immigr_c**.

7. Wage regressions

Goal

We have seen how employment varies by gender and family structure, and how gender wage gaps vary by education and immigration status. In this exercise, we will investigate the impact of all these variables on wages, using a multivariate regression.

Wages are generally not normally distributed. We will therefore apply a logarithmic transformation in order to create an outcome variable that is approximately normal, which is more suitable for regression modeling.

In addition to the variables we have already seen, we will also control for age, which has a strong relationship with earnings. Since the relationship between age and income is not necessarily linear, we will also add a term for age-squared.

Activity

Run regression models predicting **logwage**. If you have followed the instructions up to this point, you should not need to create any additional variables (see guidelines below).

You should run the regressions separately in each country, and within each country you should run a separate model for men and women. Use the following model:

$$\text{logwage} = f(\text{age agesq meduc heduc immigr partn ychild ochild ptime homeowner})$$

Produce a table of the six resulting models, with coefficients, standard errors, sample sizes, and r-squared values.

Question: Who receives a higher wage premium from being highly educated, men or women?

Question: When controlling for other individual characteristics, what is the relationship between immigrant status and wages?

Question: When controlling for other individual characteristics, do women with young children make more or less than women without children?

Guidelines

- Linear regression in base R is done with the `glm()` function. Since this command accepts weights, you can use it instead of the `svyglm()` function from the **survey** package, which is necessary for more complex sample designs.
- R will automatically detect categorical variables if they are coded as factors, meaning that you do not need to manually create dummy variables. You can also include mathematical operations directly in the definition of your regression formula if you enclose them in the `I()` function, so you do not need to manually create log wages or age squared. For example, the following estimates the log wage regression for men in the United States:

```
glm(I(log(hourwage))~age+I(age^2)+educ+immigr+livepartner+achildcat+ptime+homeowner, data=df, weights=df$ppopwgt, subset=sex=="Male" & dname=="us04")
```

- When performing several regression models in a single program, one strategy is to again write a loop that estimates each regression in turn, and then prints a summary of the model

```
model <- formula(I(log(hrwg))~age + I(age^2)+meduc+heduc+immigr+partn+ychild+ochild+ptime+ homeowner)

for(s in ...) for(d in datasets) {
  res <- glm(model, data = df, weight = df$ppopweight, subset = sex & dname = d)
  print(summary(res))
}
```

Here the definition of the model has been done before calling the loop, which makes the call to `glm()` itself somewhat easier to read. The summary function will output coefficients, standard errors and R-squared statistics.

Program

```
get_stack <- function(datasets, varp, varh, subset) {
# READ DATASETS
  pp <- read.LIS(paste(datasets, 'p', sep=''), labels=FALSE, vars=varp, subset=subset)
  hh <- read.LIS(paste(datasets, 'h', sep = ''), labels = FALSE, vars = varh)
  df <- merge(pp, hh, by = c("dname", "hid"))

# MAP NEW VARIABLES
df$homeowner <- ifelse(df$own %in% 100:199, 1, ifelse(df$own %in% 200:299, 0, NA))
df$child <- ifelse(df$children %in% c(140,200),0 ,ifelse(df$children==110,1,
ifelse(is.na(df$children),NA,2)))
df$ychild <- ifelse(df$child==1, 1, ifelse(is.na(df$child), NA, 0))
df$ochild <- ifelse(df$child==2, 1, ifelse(is.na(df$child), NA, 0))
df$partn <- ifelse(df$partner %in% c(100,120,200), 0, ifelse(is.na(df$partner),NA,1))
df$meduc <- heduc <- 0
df$meduc <- ifelse(df$educ==2 , 1, ifelse(is.na(df$educ), NA, 0))
df$heduc <- ifelse(df$educ==3 , 1, ifelse(is.na(df$educ), NA, 0))
df$hrwg <- df$gross1
df$hrwg <- ifelse(df$hrwg <= 0, NA, df$hrwg)
for (i in 1:length(datasets)) {
  topline <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],],
10*wNtile(hrwg,ppopwgt,0.5))
  df$hrwg <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],],
ifelse(df$hrwg>topline, topline, df$hrwg))
}
  return(df)
}

wNtile <- function(var, wgt, split) {
  x <- var[order(var)]
  y <- wgt[order(var)]
  z <- cumsum(y) / sum(y)
  cop <- rep(NA,length(split))
  for (i in 1:length(cop)) {
    cop[i] <- x[Find(function(h) z[h] > split[i], seq_along(z))]
  }
  return(cop)
}

#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
varh <- c('hid', 'dname', 'own')
varp <- c('hid', 'dname', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime',
'children', 'partner', 'status1', 'gross1', 'educ', 'immigr')
subset <- 'age >= 25 & age <= 54 & relation <= 2200'
df <- get_stack(datasets, varp, varh, subset)

# REGRESSION MODEL
gender <- c('Male', 'Female')
model <- formula(I(log(hrwg))~age + I(age^2) + meduc+heduc+immigr+partn+ychild+ochild+
ptime+homeowner)
for(s in (1:length(unique(df$sex[!is.na(df$sex)])))) for(d in datasets) {
  res <- glm(model, data = df, weights = df$ppopwgt, subset = sex == s & dname == d)
  print('-----')
  print(paste(gender[s], d, sep = " : "))
  print(summary(res))
}
```

Results

[1] "-----"

[1] "Male : us04"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.144e+00	8.409e-02	13.603	< 2e-16	***
age	4.854e-02	4.347e-03	11.166	< 2e-16	***
I(age^2)	-4.853e-04	5.435e-05	-8.929	< 2e-16	***
meduc	2.790e-01	1.235e-02	22.590	< 2e-16	***
heduc	6.605e-01	1.245e-02	53.040	< 2e-16	***
immigr	-6.060e-02	9.198e-03	-6.588	4.53e-11	***
partn	4.106e-02	9.720e-03	4.224	2.41e-05	***
ychild	5.658e-02	9.763e-03	5.795	6.90e-09	***
ochild	6.317e-02	8.777e-03	7.198	6.27e-13	***
ptime	-3.060e-01	1.799e-02	-17.010	< 2e-16	***
homeowner	2.411e-01	8.498e-03	28.377	< 2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 17260992 on 29985 degrees of freedom

[1] "-----"

[1] "Male : de04"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.3051407	0.2040528	6.396	1.82e-10	***
age	0.0510821	0.0104844	4.872	1.16e-06	***
I(age^2)	-0.0005623	0.0001299	-4.329	1.54e-05	***
meduc	0.0731985	0.0259923	2.816	0.00489	**
heduc	0.3297840	0.0276621	11.922	< 2e-16	***
immigr	-0.0510355	0.0234008	-2.181	0.02926	*
partn	0.0265695	0.0216033	1.230	0.21883	.
ychild	0.0449892	0.0251256	1.791	0.07345	.
ochild	0.0613726	0.0211842	2.897	0.00379	**
ptime	-0.2956306	0.0405161	-7.297	3.67e-13	***
homeowner	0.1240835	0.0174341	7.117	1.34e-12	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 2664502 on 3327 degrees of freedom

[1] "-----"

[1] "Male : gr04"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	0.7472140	0.3183373	2.347	0.019084	*
age	0.0361673	0.0163318	2.215	0.026990	*
I(age^2)	-0.0002607	0.0002006	-1.300	0.193945	.
meduc	0.1601416	0.0254636	6.289	4.55e-10	***
heduc	0.4772361	0.0276774	17.243	< 2e-16	***
immigr	-0.2503491	0.0325947	-7.681	3.40e-14	***
partn	-0.0260568	0.0377164	-0.691	0.489794	.
ychild	0.1295732	0.0331776	3.905	9.96e-05	***
ochild	0.1065245	0.0300084	3.550	0.000401	***
ptime	-0.2670276	0.0679677	-3.929	9.05e-05	***
homeowner	0.0344356	0.0235584	1.462	0.144097	.

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 176310 on 1146 degrees of freedom

[1] "-----"

[1] "Female : us04"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.124e+00	8.909e-02	12.617	< 2e-16 ***
age	4.156e-02	4.617e-03	9.003	< 2e-16 ***
I(age^2)	-4.384e-04	5.781e-05	-7.584	3.46e-14 ***
meduc	3.208e-01	1.515e-02	21.180	< 2e-16 ***
heduc	7.489e-01	1.521e-02	49.255	< 2e-16 ***
immigr	-1.975e-02	1.050e-02	-1.881	0.059924 .
partn	-1.027e-02	8.453e-03	-1.215	0.224343
ychild	3.743e-02	1.083e-02	3.456	0.000549 ***
ochild	-4.907e-02	8.510e-03	-5.766	8.20e-09 ***
ptime	-1.759e-01	9.290e-03	-18.932	< 2e-16 ***
homeowner	1.842e-01	9.112e-03	20.218	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 16539528 on 28596 degrees of freedom

[1] "-----"

[1] "Female : de04"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.9029899	0.2274761	8.366	< 2e-16 ***
age	0.0171342	0.0117989	1.452	0.14655
I(age^2)	-0.0001565	0.0001482	-1.056	0.29084
meduc	0.0873187	0.0281301	3.104	0.00192 **
heduc	0.4030372	0.0311114	12.955	< 2e-16 ***
immigr	-0.1111426	0.0248912	-4.465	8.27e-06 ***
partn	-0.0208213	0.0209912	-0.992	0.32132
ychild	0.0019161	0.0326123	0.059	0.95315
ochild	-0.0676731	0.0220161	-3.074	0.00213 **
ptime	-0.0981062	0.0195227	-5.025	5.30e-07 ***
homeowner	0.0585684	0.0191126	3.064	0.00220 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 2721877 on 3307 degrees of freedom

[1] "-----"

[1] "Female : gr04"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.2127914	0.3506948	-0.607	0.544148
age	0.0692781	0.0184864	3.748	0.000189 ***
I(age^2)	-0.0006251	0.0002336	-2.676	0.007581 **
meduc	0.3403471	0.0345974	9.837	< 2e-16 ***
heduc	0.7838179	0.0360847	21.722	< 2e-16 ***
immigr	-0.2208056	0.0421781	-5.235	2.03e-07 ***
partn	-0.0499962	0.0360101	-1.388	0.165342
ychild	0.0866716	0.0389993	2.222	0.026491 *
ochild	0.0480399	0.0339987	1.413	0.157985
ptime	-0.0991978	0.0386726	-2.565	0.010468 *
homeowner	-0.0292736	0.0282481	-1.036	0.300324

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 205048 on 962 degrees of freedom

Question: Who receives a higher wage premium from being highly educated, men or women?

- In all three countries, the coefficient for high education is higher for women, indicating a larger wage premium from having high educational attainment.

Question: When controlling for other individual characteristics, what is the relationship between immigrant status and wages?

- In general, the association between immigrant status and wages is negative, although the 95% confidence interval for the coefficient crosses 0 for German men. In general, the negative association appears stronger in Germany than in the United States, and stronger in Greece than Germany.

Question: When controlling for other individual characteristics, do women with young children make more or less than women without children?

- In general, women with young children have higher wages than women without children in the United States and Greece, but there is no association in Germany. Higher wages for women with young children could be due to a selection effect, where mothers of young children are more likely to enter the labor market if they have higher earning power.

Comments

- As we have seen, employment rates, particularly among women, vary substantially across countries. The wage regressions shown here do not account for this differential selection into employment. For this reason, many studies of wages apply a technique such as a Heckman correction, which attempts to correct for this selection bias.

8. Pooled regressions and normalised weights

Goal

In the previous exercise, we ran parallel, separate regressions for each country. In this exercise, we see an alternative approach, in which all countries are pooled together in a single model. We will continue to use classical OLS regression, but the approach shown here can easily be extended for more complex multilevel estimation approaches.

The income variables in these datasets use different currencies. To compare them, we need to convert them to a common scale. We will apply *Purchasing Power Parity* (PPP) deflators, which are intended to ensure that equal quantities of income correspond to equivalent purchasing power across currencies and national economies. The PPP deflators provided here are taken from the Penn World Tables (http://pwt.econ.upenn.edu/php_site/pwt_index.php), which are a commonly used source.

Up to this point, we have been using the weight variable **ppopwgt**, which inflates to the total population. If we use this weight in a pooled regression, every household will receive equal weight. However, this would mean that Greece — which has a much smaller population than the United States or Germany — will not have much influence on the results. In order to give each country equal weight in the model, we will use the alternative normalized weight variable **pwgt**, which always sums to 10,000 within each dataset.

Activity

Adjust wages for Purchasing Power Parities by dividing **hrwage** by the following deflators before taking the natural logarithm:

dataset	PPP deflator
US04	1
DE04	0.88
GR04	0.71

Estimate the following model, for men and women separately:

logwage = f (age agesq meduc heduc immigr partn ychild ochild ptime homeowner germany greece)

The model is the same as in the last exercise, except that it includes an indicator for country. This time, however, make sure to use *normalised*, not inflated weights.

Produce a table of the two resulting models, with coefficients, standard errors, sample sizes, and r-squared values.

Question: How can you interpret the meaning of the coefficients for the Germany and Greece dummy variables?

Question: In this pooled model, which carries a higher wage penalty: being an immigrant, or working part time?

Guidelines

- Run your regressions as you did in the previous exercise. This time, only two models need to be produced, so you may not want to use a loop.

Program

```
get_stack <- function(datasets, varp, varh, subset) {
  # READ DATASETS
  pp <- read.LIS(paste(datasets, 'p', sep = ''), labels = FALSE, vars = varp, subset =
subset)
  hh <- read.LIS(paste(datasets, 'h', sep = ''), labels = FALSE, vars = varh)
  df <- merge(pp, hh, by = c("dname", "hid"))
  # MAP NEW VARIABLES
  df$homeowner <- ifelse(df$own %in% 100:199, 1, ifelse(df$own %in% 200:299, 0, NA))
  df$child <- ifelse(df$children %in% c(140, 200) , 0, ifelse(df$children == 110 , 1,
ifelse(is.na(df$children), NA, 2)))
  df$ychild <- ifelse(df$child == 1 , 1, ifelse(is.na(df$child) , NA, 0))
  df$ochild <- ifelse(df$child == 2 , 1, ifelse(is.na(df$child) , NA, 0))
  df$partn <- ifelse(df$partner %in% c(100, 120, 200), 0, ifelse(is.na(df$partner) , NA, 1))
  df$meduc <- heduc <- 0
  df$meduc <- ifelse(df$educ == 2 , 1, ifelse(is.na(df$educ) , NA, 0))
  df$heduc <- ifelse(df$educ == 3 , 1, ifelse(is.na(df$educ) , NA, 0))
  df$ppp <- ifelse(df$dname == 'de04', 0.88, ifelse(df$dname == 'gr04', 0.71, 1))
  df$hrwg <- df$gross1
  df$hrwg <- ifelse(df$hrwg <= 0, NA, df$hrwg)
  df$germany <- ifelse (df$dname == 'de04', 1, 0)
  df$greece <- ifelse (df$dname == 'gr04', 1, 0)
  for (i in 1:length(datasets)) {
    df$hrwg <- df$hrwg / df$ppp
    topline <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],], 10 * wNtile(hrwg,ppopwgt,0.5))
    df$hrwg <- with(df[!is.na(df$hrwg) & df$dname==datasets[i],],
ifelse(df$hrwg>topline,topline,df$hrwg))
  }
  return(df)
}
wNtile <- function(var, wgt, split) {
  x <- var[order(var)]
  y <- wgt[order(var)]
  z <- cumsum(y) / sum(y)
  cop <- rep(NA,length(split))
  for (i in 1:length(cop)) {
    cop[i] <- x[Find(function(h) z[h] > split[i], seq_along(z))]
  }
  return(cop)
}
#----- RUN SCRIPTS -----
datasets <- c('us04', 'de04', 'gr04')
varh <- c('hid', 'dname', 'own')
varp <- c('hid', 'dname', 'pwgt', 'ppopwgt', 'age', 'sex', 'relation', 'emp', 'ptime', 'children',
'partner', 'status1', 'gross1', 'educ', 'immigr')
subset <- 'age >= 25 & age <= 54 & relation <= 2200'
df <- get_stack(datasets, varp, varh, subset)

# REGRESSION MODEL
gender <- c('Male', 'Female')
model <- formula(I(log(hrwg))~age + I(age^2) + meduc + heduc + immigr + partn + ychild + ochild +
ptime + homeowner + germany + greece)
for(s in (1:length(unique(df$sex[!is.na(df$sex)])))) {
  res <- glm(model, data = df, weights = df$pwgt, subset = sex == s)
  print('-----')
  print(gender[s])
  print(summary(res))
}
```

Results

[1] "-----"

[1] "Male"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.373e+00	7.104e-02	19.325	< 2e-16	***
age	4.770e-02	3.672e-03	12.990	< 2e-16	***
I(age^2)	-4.737e-04	4.561e-05	-10.387	< 2e-16	***
meduc	1.654e-01	8.218e-03	20.130	< 2e-16	***
heduc	4.969e-01	8.618e-03	57.653	< 2e-16	***
immigr	-1.092e-01	7.775e-03	-14.046	< 2e-16	***
partn	3.549e-02	7.984e-03	4.445	8.81e-06	***
ychild	6.921e-02	8.243e-03	8.397	< 2e-16	***
ochild	7.038e-02	7.294e-03	9.649	< 2e-16	***
ptime	-2.886e-01	1.485e-02	-19.438	< 2e-16	***
homeowner	1.448e-01	6.278e-03	23.059	< 2e-16	***
germany	1.805e-01	6.541e-03	27.601	< 2e-16	***
greece	1.808e-01	7.297e-03	24.781	< 2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 1146.75 on 34460 degrees of freedom

[1] "-----"

[1] "Female"

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.2524580	0.0758874	16.504	< 2e-16	***
age	0.0418405	0.0039645	10.554	< 2e-16	***
I(age^2)	-0.0004196	0.0000498	-8.425	< 2e-16	***
meduc	0.2294726	0.0097592	23.514	< 2e-16	***
heduc	0.6341410	0.0102170	62.067	< 2e-16	***
immigr	-0.0895185	0.0086977	-10.292	< 2e-16	***
partn	-0.0075337	0.0072230	-1.043	0.296953	
ychild	0.0347733	0.0093772	3.708	0.000209	***
ochild	-0.0364911	0.0072883	-5.007	5.56e-07	***
ptime	-0.1312862	0.0072164	-18.193	< 2e-16	***
homeowner	0.0828622	0.0068173	12.155	< 2e-16	***
germany	0.3098494	0.0073357	42.239	< 2e-16	***
greece	0.3346662	0.0079582	42.053	< 2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 1143.38 on 32867 degrees of freedom

Question: How can you interpret the meaning of the coefficients for the Germany and Greece dummy variables?

- These coefficients represent the overall national level of PPP-adjusted wages, when controlling for the other variables. The small negative value for Germany reflects the fact that Germany is a slightly poorer country than the United States, while the large negative value for Greece reflects the fact that Greece is substantially poorer.

Question: In this pooled model, which carries a higher wage penalty: being an immigrant, or working part time?

- The wage penalty for working part time is larger than that for being an immigrant. This is especially the case among men, although it should be noted that part time work is very uncommon among prime age men.