

# **LIS Self Teaching Package 2020**

**R version**

## **Part I**

### **Inequality, poverty, and social policy**





# Part I: Inequality, poverty, and social policy

## **Overall Plan and Structure of the Exercise**

The next eight exercises demonstrate the use of the LIS data. These exercises will lead you through the process of developing a comparative research project that examines inequality and poverty across countries. Each of the exercises introduces new concepts related to the datasets and the programming techniques needed to make use of the data. By the end of the last exercise, you will produce a complete program that returns results on poverty and inequality for a selection of five LIS countries.

Each exercise builds on the one that comes before it. It is intended that you will begin every new activity by returning to the program you have written in the previous exercise, and modifying it to satisfy the requirements of the new exercise. Each exercise contains questions that you can answer with the new results you produce. The solutions available for each exercise include an example program, with **bolded** sections indicating the code that has been added for that exercise.

The analysis shown here is simplified somewhat, compared to what might be done in an actual LIS Working Paper. Some choices have also been made in order to demonstrate particular aspects of the LIS data. However, these exercises provide a starting point for researchers who want to develop an analysis of the data based on their own research questions.

## **Research Questions**

Since the beginning of the LIS project, one of the most prominent objects of research using the data has been the effect of government tax and transfer programs on poverty and inequality. The first substantive paper in the LIS Working Papers series, published in 1985, analysed the pre- and post-transfer poverty rate in Sweden, the United Kingdom, Israel, the United States, Norway, Canada, and West Germany.<sup>1</sup> The second substantive paper compares the distribution of income across these same seven countries.<sup>2</sup>

---

1 Richard Hauser, Lee Rainwater, Martin Rein, Gaston Schaber, Timothy Smeeding. "Poverty in Major Industrialized Countries". LIS Working Paper No. 2 – Jul 1985. Available at: <http://www.lisdatacenter.org/wps/liswps/2.pdf>

2 Michael O'Higgins, Gunther Schmaus, Geoffrey Stephenson. "Income Distribution and Redistribution". LIS Working Paper No. 3 – Jun 1985. Available at: <http://www.lisdatacenter.org/wps/liswps/3.pdf>

States affect the income distribution through several different types of policy, which are captured in the LIS data:

- Income taxes and social insurance contributions
- Social insurance programmes linked to employment, such as public pensions, unemployment insurance benefits, and sickness pay
- Universal benefits provided irrespective of employment, income or assets
- Social assistance benefits for especially needy individuals or households

The exercises in this section assess the impact of these policies in different countries on both poverty and income inequality. We will measure income from labour, capital, private pensions, and private transfers. Then, we will compare them to income after income taxes and social insurance contributions as well as government transfers are accounted for. Within the category of government policies, we will separate the effect of payroll taxes, social insurance, and universal benefits from the effect of social assistance benefits. We will measure the proportion of the population that is poor according to these different income measures, in which poverty is defined relative to median level of income within a country. We will also compare income inequality using one of the most popular and longstanding inequality measures, the Gini coefficient.

As LIS has grown, the analysis of government policy, poverty, and inequality has been updated for more countries and more recent years.<sup>3</sup> Recently, LIS expanded beyond the rich countries that have long made up the core of the project, and began adding data from middle income countries. This gives researchers the opportunity to compare income and poverty in these countries to the patterns seen in rich countries that have been much more heavily studied.

For this exercise, we will begin by analysing the data from Guatemala. We will then compare Guatemala to four other countries: The United States, Denmark, Hungary, and Israel. These countries have been chosen in part for pedagogical reasons. As well as, they allow for substantively interesting comparisons as they represent a wide range of national income levels and welfare state regime types.

After completing the final exercise, you will be able to answer the following questions:

- Which of the five countries in the study have the highest levels of income

---

<sup>3</sup> For a recent example, see Timothy Smeeding, "Government Programs and Social Outcomes: The United States in Comparative Perspective". LIS Working Paper No. 426 – May 2005. Available at: <http://www.lisdatacenter.org/wps/liswps/426.pdf>

inequality and poverty before taxes and transfers are accounted for? Does this change when taxes and transfers are included?

- Which type of government policy has a larger impact on inequality and poverty in each country: taxes, social insurance, and universal benefits, or targeted social assistance?

### **Before you begin**

Before you begin the exercises, take a look at the [2019 Template LIS User Guide](#), which can be accessed through *LIS Website* → *Our Data* → *LIS Database*. The User Guide provides an overview of the structure of LIS data and some data management practices, such as missing values policy and aggregation rules, which will be useful in working with LISSY.

In addition to this, an overview of the datasets and variables is provided through the METadata Information System ([METIS](#)) without having to login to LISSY. You can access METIS via the *LIS Website* → *METIS* → *Enter METIS* → *LIS*. After selecting the datasets and variables, consult the *Results tab* for information on variables and definitions, dataset-specific information and variable availability across datasets.

## Contents

### **1. Accessing the LIS databases LISSY Interface**

- Submit a LIS job and get basic descriptive statistics for one variable

### **2. Sample selection and weighting**

- Introduce the list of variables to be used
- Select the sample and eliminate cases with missing data
- Produce weighted and unweighted descriptive statistics for variables

### **3. Working with household income variables: top and bottom coding and equivalence scales**

- Top and bottom coding income data to remove outliers
- Correcting income for household size using equivalence scales

### **4. Inequality: The Gini Index**

- Calculating the Gini coefficient

### **5. Relative poverty rates**

- Calculating relative poverty

### **6. Comparing income concepts**

- Introducing four concepts of income
- Inequality and poverty before and after taxes and transfers

### **7. Comparing multiple countries**

- Extending the analysis to multiple countries
- Working with net vs. gross income data

### **8. Producing compact and concise output**

- Produce easy-to-use output

### **9. Producing graphs with LISSY**

- Generate Lorenz curves for multiple countries

# 1. Accessing the LIS Database: LISSY Interface

## Goal

This exercise introduces [LISSY Interface](#), which we will be using to work with LIS data in all of the subsequent exercises.

LISSY secure web-based interface allows researchers to:

- write, submit and view job requests (and corresponding outputs);
- track the status of the job requests in process ('received', 'processing', 'set for review', 'refused', etc.); and
- access the history of all job requests ever sent.

In this exercise we will use [LISSY interface](#) to open a dataset and produce basic descriptive statistics.

## Activity

Go to *LIS website* → [Login LISSY](#) tab with your LISSY account. Submit a simple program to display descriptive statistics (number of valid observations, mean, minimum, maximum) for the household-level income variable **dhi**, for Guatemala 2006. **dhi**, or disposable household income, contains the total monetary and non-monetary current income for the household, net of income taxes and social security contributions. It is a harmonised variable that is available for all datasets.

## *Questions*

1.1. How many valid observations (non-missing) are in this dataset for **dhi**?

## Guidelines

- Once connected to LISSY Interface, there are three main tasks that may be carried out:
  1. Submit jobs through the Job Session window >>edit job pane.
    - Select a project (LIS, LWS).
    - Select a statistical package (SAS, SPSS, Stata or **R**).
    - When submitting a job (edit job pane), always add a subject line.
    - Write your code.
    - Click on the submit button.
  2. Work with Today's Jobs (Recent Jobs window.)
    - Watch the status of jobs currently sent to LISSY in the 'recent jobs'

pane.

- View the jobs returned by LISSY.
- Click on the received job (marked with green arrow).
- Click on the 'job text' or 'listing' tabs, respectively, of the right panel to see the request and its output.
- Download the received results and the job code by either:
  - clicking on 'save all results to pdf file' button in the middle of the upper tool bar to save the results in pdf format.
  - Clicking on 'save all results to a zipped file' button to save the results in txt format or png format for the produced graphs.
- Re-submit a selected job by clicking on the 'upload all or selected text to the actual job editor' button at the left hand side of the upper tool bar.

### 3. Manage (view, clean and search) all job requests ever sent in the Job Archive window.

- View jobs sent over a specific time period.
- Clean the library by discarding useless job requests (remove job from active archive button).
- Search jobs by keywords.
- Re-submit a selected job by clicking on the 'upload all or selected text to the actual job editor' button at the left hand side of the upper tool bar.

- When you open a LIS dataset, use the correct file reference for the country/year you wish to use. For example:

```
df ← read.LIS('gt06p')
```

For more information about the syntax of country/year file reference, see the job submission instructions on the LIS web site (*Data Access → Job Submission*), accessed through this [link](#).

For a list of available data sets and their 2digit country codes, go to:

*Our Data → LIS Database → List of Datasets*, accessible through [LIS](#)

*Our Data → LWS Database → List of Datasets*, accessible through [LWS](#)

- R reminders

- To access a variable within a data frame, you can use either the syntax **df\$<variable>** or **df[['<variable>']]**
- You can get basic descriptive statistics for a variable (such as the



minimum, maximum, mean and median) using `summary(<variable>)`. The option `digits` is needed to ensure that the display precision is great enough, `summary(<variable>, digits=<#>)`.

- To get the total number of non-missing observations of a variable, you can use `sum(!is.na(<variable>))`. This creates a vector which is TRUE whenever the variable is not missing, and then gives you the sum of that vector, i.e., the number of non-missing observations.
- If you do not receive your job in the expected amount of time, it means that there is a long queue of jobs on LISSY. In that case, resending your job, or sending several other ones while waiting to receive the first one, will only increase the queue and hence your waiting time. Remember to wait to get your results before sending a new job!

## Program

```
df <- read.LIS('gt06h')
print(summary(df$dhi, digits=10))
print(sum(!is.na(df$dhi)))
```

## Results

	Number of valid observations	Mean	Minimum	Maximum
dhi	13,664	40,670	-186,010	2,729,298

### *Solution*

- 1.1. How many valid observations (non-missing) are in this dataset for **dhi**?
  - There are 13,664 valid observations.

### *Comments*

- It is important to pay close attention to sample sizes in order to make sure you have enough data to make stable estimates. When working with small datasets, or small sub-samples of datasets, always check the size of the sample underlying each statistic you have computed.
- As you can see in the results, the disposable household income **dhi** can be negative. This happens in cases where the data provider kept losses as such rather than applying bottom coding techniques. This typically happens with incomes from self-employment or capital income; in rare circumstances it happens that taxes are higher than gross income due to different income reference periods or miscalculation of taxes.

## 2. Sample selection and weighting

### Goal

This exercise introduces the variables that we will be using in the rest of our analysis and concentrates on sample selection and the use of weights.

*Sample selection* - The final objective of this exercise is to compare incomes before and after government intervention. In order to be sure that the comparison is correct, users should ensure that they use exactly the same sample when calculating statistics for the pre- and post-government intervention. It is thus important to begin by selecting a sample which can be used for the entire analysis. For this reason, we will drop from the analysis not only cases which have missing values in the variable of interest to the specific statistic being calculated at each step, but all the cases that have a missing value in any of the variables of interest for the whole analysis.

*Use of weights* - Comparative researchers are typically interested in the characteristics of national populations, not the samples provided. It is very important to understand and use sample weights correctly in order to get representative results for the total underlying population. This exercise shows the differences in statistics between the unweighted sample and the weighted population.

### Activity

As in the previous exercise, we will continue working with the Guatemala 2006 (GT06) data. Modify your previous program to select only the following variables in addition to **dhi**: household weight (**hpopwgt**), number of household members (**nhhmem**), gross or net income information (**grossnet**), factor income (**hifactor**), work-related insurance transfers (**hpub\_i**), universal benefits (**hpub\_u**), social assistance benefits (**hpub\_a**), private transfers (**hiprivate**) and tax and social security contribution expenditures (**hxitsc**).

(1) Produce two different sets of descriptive statistics for the variables you have selected:

- For continuous variables, the statistics - unweighted, and weighted by person - should include the number of observations, the mean, the median, the minimum and the maximum;
- For categorical variables, you should produce a frequency table.

(2) Drop all cases for which **dhi** or any of its income and expenditure sub-components is missing and see how many cases have missing data.

## Questions

- 2.1. What currency unit is used for the income variables in this dataset?
- 2.2. What effect does applying the weights have on the median of disposable household income?
- 2.3. What percentage of cases is being dropped from the dataset?
- 2.4. In subsequent exercises, we will compare incomes before and after accounting for different kinds of government transfers. What is the difference, in terms of definition, between work-related insurance transfers, universal benefits, and assistance benefits?
- 2.5. In subsequent exercises, we will separate social assistance transfers from work-related insurance and from universal transfers. Based on the information you have so far, which type of transfer do you think has a larger effect on inequality and poverty measures in Guatemala?
- 2.6. How many different values does the variable **grossnet** take in this dataset? How might this variable be useful?

## Guidelines

- The income concepts used for LIS income variables are defined in the Interactive METadata Information System (METIS).
- Note to R users: R is a very flexible language, and there are typically many ways to code the same operation. The tips in these lessons suggest one method of achieving some of the exercises goals, but if you are familiar with other techniques you should feel free to use them.
  - *Note that for these early lessons, which concentrate on simple descriptive statistics, R may be somewhat more cumbersome to use than other statistical packages which are designed to make such simple estimates very easy to retrieve. This is in contrast to the later lessons, where R's power and flexibility make it possible to produce more complex estimates very quickly and compactly.*
- Use of weights in basic descriptive is not included in R base but is a part of R packages, such as **Hmisc**,

```
wmean <- function(x, weight) {  
  y <- x[which(!is.na(x))]  
  wgt <- weight[which(!is.na(x))]  
  wmean <- sum(y*wgt/sum(wgt))  
  return(wmean)  
}
```

- *Note that this function is not robust in the sense that it does not test for missing arguments when it is called and ... so forth. But you may want to re-use and improve it later for your own coding purpose*
- Use the option "labels=FALSE" with **read.LIS** to return only numeric codes rather than value labels, which will make recoding easier later on.
- You can add the "vars" argument to the **read.LIS** function to keep only certain variables (columns):  
**df <- read.LIS('gt06', labels=FALSE, vars=<varlist>)**
- This avoids placing unnecessary burden on the machine so that the submitted jobs will run faster. You can also use the "subset" to keep only certain observations. For example, to select only cases with non-missing household income:

**df <- read.LIS('gt06', labels=FALSE, subset="complete.cases(dhi)")**

## Program

```
wmean <- function(x, weight) {
  y <- x[which(!is.na(x))]
  wgt <- weight[which(!is.na(x))]
  wmean <- sum(y*wgt/sum(wgt))
  return(wmean)
}

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)
}

vars <- c('dhi','hifactor','hpub_i',
'hpub_u','hpub_a','hiprivate','hxitsc','hpopwgt','nhhmem','grossnet')
df <- read.LIS('gt06h', labels=FALSE, vars=vars)
subset <- 'complete.cases(dhi,hifactor,hpub_i,hpub_u,hpub_a,hiprivate,hxitsc)'
df2 <- read.LIS('gt06h', labels=FALSE, vars=vars, subset=subset)
print(row_total <- nrow(df))
print(row_drop <- with(df,
length(which((complete.cases(dhi,hifactor,hpub_i,hpub_u,hpub_a,hiprivate,hxitsc)
== TRUE))))))
print(miss_income<- row_total- row_drop )
round(((row_total - row_drop) / row_total) * 100, digits = 2)
print(c(toupper('grossnet')))
print(table(df$grossnet, useNA = 'ifany'))
print(paste(round(prop.table(table(df$grossnet, useNA = 'ifany')) * 100, digits =
2), "%", sep = ""))
cat(paste(" "), sep = '\n')
for (x in c('hpopwgt','dhi','hifactor','hpub_i', 'hpub_u','hpub_a','hiprivate','hxitsc')) {
df1 <- df[!is.na(df[[x]]), ]
print(c(toupper(x)))
print(c(nb_obs = sum(!is.na(df1[[x]]))))
print(summary(df1[,x], digits = 10))
```

```

print(c(weighted_mean = round(wmean(df1[[x]], df1$hpopwgt*df1$nhhmem),
digits = 2), weighted_median = round(wNtile(df1[[x]], df1$hpopwgt *
df1$nhhmem, split = 0.5), digits = 2)))
  cat(" ", sep = '\n')
}
for (x in c('hpopwgt','dhi','hifactor','hpub_i', 'hpub_u','hpub_a','hiprivate','hxitsc'))
{
print(c(toupper(x)))
print(c(nb_obs = sum(!is.na(df2[[x]]))))
print(summary(df2[,x], digits = 10))
print(c(weighted_mean = round(wmean(df2[[x]], df2$hpopwgt*df2$nhhmem),
digits = 2), weighted_median = round(wNtile(df2[[x]], df2$hpopwgt *
df2$nhhmem, split = 0.5), digits = 2)))
  cat(" ", sep = '\n')
}

```

## Results

	number of observations	percent
Valid	13,664	99.84
Missing	22	0.162
Total	13,686	100

grossnet	number of observations	percent	cumulative percent
[120]gross, taxes and contributions imputed	13,686	100	100
Total	13,686	100	

### Unweighted, all cases

	number of observations	mean	median	minimum	maximum
hpopwgt	13,686	193.8	119	2	2,657
dhi	13,664	40670	27,092	-186,010	2,729,298
hifactor	13,664	39,924	24,778	0	3,866,644
hpub_i	13,664	803.5	0	0	134,400
hpub_u	13,686	0	0	0	0
hpub_a	13,686	695.5	0	0	62,400
hiprivate	13,664	1,439	0	0	128,020
hxitsc	13,664	2,192	0	0	1,137,346

### Weighted, all cases

	mean	median
hpopwgt	443.94	295
dhi	49,556.30	33658
hifactor	49,298.82	30819
hpub_i	777.41	0
hpub_u	0	0
hpub_a	757.02	0
hiprivate	1,675.62	0
hxitsc	2,952.58	0

## Weighted, missing income cases dropped

	mean	median
hpopwgt	443.705	295
dhi	49,556.30	33,658
hifactor	49,298.82	30,819
hpub_i	777.41	0
hpub_u	0	0
hpub_a	757.02	0
hiprivate	1,675.62	0
hxitsc	2,952.58	0

### Solution

2.1. What currency unit is used for the income variables in this dataset?

- The currency unit in this dataset is the Guatemalan Quetzal. This information can be found on the LIS METIS documentation system at *Our Data* → *METIS*, accessed through this [link](#). One way to find currency information in [METIS](#) is by *clicking on enter Metis* → *selecting LIS* → *selecting a country* → *Results* → *Dataset information* → *[LIS] Dataset*.
- Each LIS dataset has a variable named "currency". This variable indicates the currency unit used for the income/expenditure variables. This variable is available in both the household and the person data files.

```
df <- read.LIS('gt06h')
print(summary(df$currency, digits=10))
```

currency units	number of observations	Percent	Cum.
[320]GTQ - Quetzal	13,664	100	100
Total	13,664	100	

2.2. What effect does applying the weights have on the median of disposable household income?

- The median of unweighted **dhi** is 27,092 Quetzals, while the median of weighted **dhi** is 33,658 Quetzals, suggesting that low-income households are over-represented in the sample.

2.3. What percentage of cases is being dropped from the dataset?

- There are 22 cases with missing data, or 0.16 percent of the total number of cases in the dataset. Note that you should always be careful if you see a large amount of missing data, as it could bias your estimates.

2.4. In subsequent exercises, we will be comparing incomes before and after



accounting for different kinds of government transfers. What is the difference between work-related insurance transfers, universal benefits, and assistance benefits?

- These concepts are defined in the [Interactive METadata Information System \(METIS\)](#).

*LIS Database* → *Select variables* → *Major economic aggregates* → *hpub\_i/hpub\_u/hpub\_a* → *click on the information icon ⓘ* → a pop-up menu will appear displaying all the information related to these variables (definition and comments):

- Work-related insurance transfers: Monetary transfers stemming from systems where the eligibility is based on the existence and/or the length of an employment relationship; in most cases the benefits are financed by contributions paid by employers, workers or both, and their amount is usually dependent on either the previous earnings or the previous contributions (not including occupational and voluntary pensions);
- Universal transfers: Monetary transfers stemming from public programmes that provide flat-rate benefits to certain residents or citizens, provided that they are in a certain situation, but without consideration of income, employment or assets; note that in some cases the benefit amount may also depend on the other incomes of the individuals, which at the limit may result on some proportion of the population at the upper end of the income distribution to be excluded from receipt;
- Assistance transfers: Monetary and non-monetary transfers stemming from public programmes that provide benefits especially targeted to needy individuals or households (i.e. with a strict income or assets test); the amount of the benefits is either flat rate or based on the difference between the recipient income and a standard amount representing the minimum subsistence needs as guaranteed by the government.

2.5. In subsequent exercises, we will separate social assistance transfers from work-related insurance and from universal transfers. Based on the information you have so far, which type of transfer do you think has a larger effect on inequality and poverty in Guatemala?

- Since universal transfers are not available at the level of detail required for Guatemala's dataset. Therefore, inequality and poverty will only be impacted by social assistance transfers and work-related insurance.

2.6. How many different values does the variable **grossnet** take in this

dataset? How might this variable be useful?

- The variable **grossnet** has the same value for every case ([120] gross, taxes and contributions imputed). By looking at this variable, you can learn what type of gross or net income information is contained in this variable, even without looking at the documentation. This variable will also be useful when working with multiple datasets that may contain either gross or net income. This information can be found in [METIS](#) under the *Results* → *Dataset information* → *Code Books* tab.

### 3. Working with household income variables: top and bottom coding and equivalence scales

#### Goal

In order to compare incomes across countries, we need to make sure that our variables are fully comparable. In this exercise, you will apply top- and bottom-codes to remove extreme values. You will then create an *equivalised* income variable that adjusts for household size.

*Top and bottom coding* - Many inequality measures are sensitive to the values at the bottom and/or top of the income distribution, and some are not defined for non-positive values of income (e.g., any measure that calculates a logarithm). Applying top and bottom-codes (often referred to as 'winsorising') will avoid this problem, as well as ensuring comparability between datasets that may have originally had different top- and bottom-coding.

*Equivalence scales* - In order to get measures of poverty and/or income inequality in a population, it is necessary to compare income across different types of households. It is not logical to directly compare total household income between households of different sizes and composition.

Suppose you observe three levels of income (A, B, and C), where  $A > B > C$ . You cannot state that a household earning A is better off than one earning B unless you know the two households are similar in composition. For example, a family of four adult members receiving A is not necessarily better off than a couple with two children who receive B, and the family receiving B may not be better off than the childless couple receiving C.

For this reason, total household income needs to be adjusted to make it comparable across different households. This exercise gives one example of "equalising" households using one specific equivalence scale.

#### Activity

- (1) Keep the code from your previous exercise that you used to drop cases with missing data.
- (2) Create a new variable, **dhiT**, a top- and bottom-coded version of **dhi**. Bottom-code by setting all values less than zero to zero. Top-code by setting all values greater than ten times the median of **dhi** to ten times the median of **dhi**.
- (3) Create another new variable, **edhi\_tb**, an equivalised version of top- and bottom-coded disposable household income. We will correct for household size by applying the "LIS equivalence scale" (i.e., the square root of the number of household members).

- (4) Create a measure of per-capita income, **cdhi**, by dividing household income (un-equivalised, but top- and bottom-coded) by the number of household members.
- (5) Produce summary statistics showing the mean, median, minimum, and maximum of the four income variables: **dhi**, **dhiT**, **cdhi\_tb**, and **edhi**.

Using the results returned from LISSY, fill out the following table:

	Household income (no top or bottom codes)	Household income	Per capita income	Equivalised income
Mean				
Median				
Minimum				
Maximum				

### Questions

- 3.1. Which of these four versions of the income variable contain negative values?
- 3.2. Relative to household income and per capita income, how large are the mean and median of equivalised income?
- 3.3. How does applying the top and bottom code affect the mean and median of household income?

### Guidelines

- As we continue to build up to our final program, some of the code from the previous exercises will no longer be necessary. (For example, the code that produced the summary statistics in the previous exercise). You can choose to delete this code from your program in order to make it shorter. However, if you would like to keep a line code but stop it from being executed, simply place a **#** before it.
- This exercise does not use all of the variables that were used in the previous section. But you should continue to keep all of those variables when you open the dataset as they will be needed in future exercises.
- To equalise income, divide the total household income by the value of the equivalence scale for each observation. To generate LIS equivalised income:
 

```
df$edhi <- df$dhiT / (df$nhhmem^0.5)
```
- Be careful when using weights. Make sure that the weight matches your unit

of analysis. Weigh by **hpopwgt** for variables which are intrinsically at the household level (e.g., **dhi**) and by **hpopwgt\*nhhmem** (to account for household size) for variables that are conceptually meaningful at the person level (e.g., per capita and equivalised income).

## **Program**

```
wmean <- function(x, weight=NULL) {
  if (is.null(weight))
    weight <- rep(1, length(x))
  y <- x[which(!is.na(x))]
  wgt <- weight[which(!is.na(x))]
  wmean <- sum(y*wgt/sum(wgt))
  return(wmean)
}

wNtile <- function(var, wgt = NULL, split) {
  if (is.null(wgt))
    wgt <- rep(1, length(var))
  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)
}

topBottom <- function(var, botline, topline) {
  tb <- ifelse(var < botline, botline, var)
  tb[tb > topline] <- topline
  return(tb)
}

setups <- function(data_file) {
  vars <- c('dhi', 'hifactor', 'hpub_i', 'hpub_u', 'hpub_a', 'hiprivate', 'hxitsc',
'hpopwgt', 'nhhmem', 'grossnet')
  subset <- 'complete.cases(dhi, hifactor, hpub_i, hpub_u, hpub_a, hiprivate,
hxitsc)'
  df <- read.LIS(data_file, labels=FALSE, vars=vars, subset=subset)
  botline <- 0
  topline <- 10 * wNtile(df$dhi, df$hpwgt, 0.5)
  df$dhiT <- topBottom(df$dhi, botline, topline)
  df$edhi <- df$dhiT / df$nhhmem^0.5
  df$cdhi <- df$dhiT / df$nhhmem
  return(df)
}
```

```

df <- setups('gt06h')
for (x in c('dhi', 'dhiT')) {
  cat(paste("VARIABLE: ", toupper(x), sep=""), sep = '\n')
  cat(paste("Average: ", format(round(wmean(df[[x]] , df$hpopwgt) , digits =
2), big.mark = ", ")), sep = '\n')
  cat(paste("Median : ", format(round(wNtile(df[[x]], df$hpopwgt, 0.5), digits = 2),
big.mark = ", ")), sep = '\n')
  cat(paste("Minimum: ", format(round(min(df[[x]]), digits = 2), big.mark = ", ")),
sep = '\n')
  cat(paste("Maximum: ", format(round(max(df[[x]]), digits = 2), big.mark = ", ")),
sep = '\n')
  cat(" ", sep = '\n')
}
for (x in c('cdhi', 'edhi')) {
  cat(paste("VARIABLE: ", toupper(x), sep=""), sep = '\n')
  cat(paste("Average: ", format(round(wmean(df[[x]] , df$hpopwgt *
df$nhhmem) , digits = 2), big.mark = ", ")), sep = '\n')
  cat(paste("Median : ", format(round(wNtile(df[[x]], df$hpopwgt * df$nhhmem,
0.5), digits = 2), big.mark = ", ")), sep = '\n')
  cat(paste("Minimum: ", format(round(min(df[[x]]), digits = 2), big.mark = ", ")),
sep = '\n')
  cat(paste("Maximum: ", format(round(max(df[[x]]), digits = 2), big.mark = ", ")),
sep = '\n')
  cat(" ", sep = '\n')
}

```

## Results

	Household income (no top or bottom codes) <b>(dhi)</b>	Household income (top or bottom coded) <b>(dhi_tb)</b>	Per capita income (top or bottom coded) <b>(pcdhi_tb)</b>	Equivalentised income <b>(edhi_tb)</b>
Mean	47,681.16	46,027.95	9,416.68	20,665.61
Median	31,096	31,096	5,766	14,007
Minimum	-186,010	0	0	0
Maximum	2,729,298	310,960	293,335	293,335

### *Solution*

- 3.1. Which of these four versions of the income variable contain negative values?
- Only the income variable without top or bottom codes contains negative values, which are removed by applying a bottom-code. This will be important in the next exercise on inequality. The measure of inequality we will be using, the Gini coefficient, does not allow negative values. Removing negative values also allows for the commonly-used logarithmic transformation of income.
- 3.2. Relative to household income and per capita income, how large are the mean and median of equivalentised income?
- The mean and median values for equivalentised income fall between those for household income and those for per capita income. The equivalentising formula of  $\mathbf{dhi}/(\mathbf{nhmem}^{0.5})$  is a compromise between assigning all individuals their household income ( $\mathbf{dhi}/\mathbf{nhmem}^0$ ) and assigning them a per capita income ( $\mathbf{dhi}/\mathbf{nhmem}^1$ ).
- 3.3. How does applying the top and bottom code affect the mean and median of household income?
- Applying top and bottom codes makes the mean lower but does not affect the median (31,096). Means can be very sensitive to extreme values, so median values are often preferred as a measure of central tendency.



## 4. Inequality: The Gini Index

### Goal

This exercise introduces the Gini index, which is one of the most commonly used income inequality indicators. We will be using the Gini coefficient as our measure of inequality in subsequent exercises, in order to compare inequality across countries and across different concepts of income.

### Activity

**Calculate the Gini index on total disposable income for Guatemala in 2006, using variables created in the previous exercise.**

- (1) Start with your program from the previous exercise, which will drop observations with missing data,
- (2) Apply top and bottom codes,
- (3) Create variables containing equivalised disposable income and disposable income per capita,
- (4) Calculate the Gini coefficient for the winsorised (or bottom- and top-coded) versions of household income, per capita income, and equivalised income, and fill out the following table:

	Household income	Per capita income	Equivalised income
Gini coefficient			

### *Questions*

- 4.1. Which measure shows greater inequality: household income, per capita income or equivalised income? What does this suggest about the possible relationship between income and household size?

### Guidelines

- To make the most of R, it is advisable to break your code into functions as much as possible. In future exercises, we will adapt our program in order to make use of functions. For this exercise, we will provide a function definition that you can include in your code to produce weighted Gini coefficients:

```
gini <- function(x,weight) {  
  ox <- order(x)  
  x <- x[ox]  
  weight <- weight[ox]/sum(weight)
```

```
p   <- cumsum(weight)
nu  <- cumsum(weight*x)
n   <- length(nu)
nu  <- nu/nu[n]
res <- sum(nu[-1]*p[-n])-sum(nu[-n]*p[-1])
return(res)
}
```

- After it has been defined, you can call this function just as you would call any built-in R function, for example:

```
gini(df$dhi, df$hpopwgt)
```

- This function is derived from the **reldist** package, <http://cran.r-project.org/web/packages/reldist/>.

## Program

```
gini <- function(x,weight) {  
  ox  <- order(x)  
  x   <- x[ox]  
  weight <- weight[ox]/sum(weight)  
  p    <- cumsum(weight)  
  nu   <- cumsum(weight*x)  
  n    <- length(nu)  
  nu   <- nu/nu[n]  
  res  <- sum(nu[-1]*p[-n])-sum(nu[-n]*p[-1])  
return(res)  
}  
  
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)  
}  
  
topBottom <- function(var, botline, topline) {  
  tb      <- ifelse(var < botline, botline, var)  
  tb[tb > topline] <- topline  
  return(tb)  
}  
  
setups <- function(data_file) {  
  vars    <- c('dhi', 'hifactor', 'hpub_i', 'hpub_u', 'hpub_a', 'hiprivate', 'hxitsc', 'hpopwgt',  
'nhhmem', 'grossnet')  
  subset  <- 'complete.cases(dhi, hifactor, hpub_i, hpub_u, hpub_a, hiprivate, hxitsc)'  
  df     <- read.LIS(data_file, labels=FALSE, vars=vars, subset=subset)  
  botline <- 0  
  topline <- 10 * wNtile(df$dhi, df$hpopwgt, 0.5)  
  df$dhi <- topBottom(df$dhi, botline, topline)  
  df$edhi <- df$dhi / df$nhhmem^0.5  
  df$cdhi <- df$dhi / df$nhhmem  
  return(df)  
}
```

```
df <- setups('gt06h')
round(gini(df$dhi, df$hpopwgt)      , digits = 4)
round(gini(df$cdhi, df$hpopwgt*df$nhhmem), digits = 4)
round(gini(df$edhi, df$hpopwgt*df$nhhmem), digits = 4)
```

## **Results**

	Household income	Per capita income	Equivalised income
Gini coefficient	0.480	0.510	0.470

### *Solution*

4.1. Which measure shows greater inequality: household income, per capita income or equivalised income? What does this suggest about the possible relationship between income and household size?

- The Gini is lower for equivalised income, than the Gini for household income, which in turn is lower than the Gini for per-capita income. This reflects the fact that poorer households in Guatemala tend to be larger than richer households. Because the LIS equivalence scale assumes some economies of scale in large households, it produces a lower estimate of inequality than a per-capita measure.

## 5. Relative poverty rates

### Goal

In order to calculate any measure of poverty, it is essential to make some assumptions concerning the criteria used to define poverty. The approach used by LIS (and most commonly adopted in the literature), is that of creating a relative poverty line based on the level and distribution of (equivalised) household disposable income in the total population. Households are classified as poor or non-poor on the basis of whether their income is lower or higher than the relative poverty line.

Once poor households are identified, you can create an indicator to help identify the proportion of poor households (or individuals) and to measure the level of poverty. The choice of the indicator used will mainly depend on the purpose of the research. In this exercise, we will calculate the relative poverty rates of households and individuals in the Guatemala 2006 data.

### Activity

- (1) Add code to your program to produce an indicator for poverty.
- (2) Define the poverty line as 50% of the median equivalised income.
- (3) Calculate both the percentage of households in poverty and the head count ratio (defined as the percentage of individuals living in poor households), and complete the following table.

	Households	Individuals
Relative poverty rate		

### *Question*

- 5.1. Are there more poor *households* or more poor *individuals*? What can you infer from this?

### Guidelines

- From this point forward, we will be working exclusively with equivalised income, so the sections of your code relating to per-capita income can now be commented out or removed.
- In R, you can take the sum of a logical vector, and the result will be the proportion of elements in that vector which have the value TRUE. This means that you can place a logical comparison inside a sum statement to produce a poverty rate, as in the following example:

```
sum((df$edhi < 0.5 * wNtile(df$edhi, df$popwgt * df$nhmem, 0.5)) *
```

**df\$hpopwgt) / sum(df\$hpopwgt))**

This code computes the weighted median income, sum the equivalised disposable income that is for those cases that are below half of that median, and then divides the result by the sum of the weights, which ends up to the percentage in relative poverty.

- Whether your poverty rate is the proportion of *individuals* or *households* in poverty depends on which weighting you use. Use **hpopwgt** if you want to measure household poverty, and **hpopwgt\*nhhmem** if you are interested in individual poverty. If you use the individual-level weighting, you will produce the Head Count Ratio (HCR), which is the percentage of poor individuals in the total population.

## **Program**

```
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)
}

topBottom <- function(var, botline, topline) {
  tb <- ifelse(var < botline, botline, var)
  tb[tb > topline] <- topline
  return(tb)
}

setups <- function(data_file) {
  vars <- c('dhi', 'hifactor', 'hpub_i', 'hpub_u', 'hpub_a', 'hiprivate', 'hxitsc', 'hpopwgt',
            'nhhmem', 'grossnet')
  subset <- 'complete.cases(dhi, hifactor, hpub_i, hpub_u, hpub_a, hiprivate, hxitsc)'
  df <- read.LIS(data_file, labels=FALSE, vars=vars, subset=subset)
  botline <- 0
  topline <- 10 * wNtile(df$dhi, df$hpowwgt, 0.5)
  df$dhi <- topBottom(df$dhi, botline, topline)
  df$edhi <- df$dhi / df$nhhmem^0.5
  df$cdhi <- df$dhi / df$nhhmem
  return(df)
}

df <- setups('gt06h')
maxline <- 0.5
round(100 * (sum((df$edhi < maxline * wNtile(df$edhi, df$hpowwgt * df$nhhmem,
0.5)) * df$hpowwgt) / sum(df$hpowwgt)), digits = 2)
round(100 * (sum((df$edhi < maxline * wNtile(df$edhi, df$hpowwgt * df$nhhmem,
0.5)) * df$hpowwgt * df$nhhmem) / sum(df$hpowwgt * df$nhhmem)), digits =
2)
```



## **Results**

	Households	Individuals
Relative poverty rate	22.1%	22.5%

### *Solution*

5.1. Are there more poor households or more poor individuals? What can you infer from this?

- There is a slightly greater proportion of poor individuals than poor households. This is because poor households are larger on average than non-poor households. This is another reason why the use of the equivalence scale is important.

### *Comments*

- The head count ratio (HCR) measures poverty incidence (i.e., the number or proportion of poor people), but gives every person equal weight no matter how far they fall from the poverty line.
- Another measure, the Income Gap Ratio (IGR) measures poverty intensity or depth (how poor are the poor), but one poor person with an income of an amount  $x$  counts the same as two poor people each with an income of  $x/2$ . That is, the IGR measures the average income gap, but not its distribution among the poor.
- There are many other indicators of poverty that may be useful for different purposes. These include, among the most common, the whole family of Foster-Greer-Thorbecke indicators (of which the HCR is only one), the Sen index, the Takayama index, the Clark index, and the Thon index. It is important to note that a country may score better in comparison to a second country when using a particular index, but could score worse if another index was used instead.

## 6. Comparing income concepts

### Goal

Now that we have calculated Gini coefficients and poverty rates based on equivalised household disposable income, we can easily apply this same code to three other income concepts: income before any taxes and government transfers, income after taxes, social insurance, and universal benefit transfers, and income after social assistance transfers. Starting from this exercise, we will also introduce some programming techniques which will make it easier to repeat a series of commands several times without having to repeat the code.

*Different income concepts* - The income variable we have been working with, **dhi**, combines multiple income and expenditure flows. It is the sum of labour and capital income, private transfers, private pensions, work-related insurance transfers, universal benefits, and social assistance benefits, minus any taxes and social insurance contributions paid. We will now define three new concepts of income. One of them is income before *any* government redistribution, but including private pensions. The second is income *after* taxes, social insurance, and universal benefits, but *before* social assistance is included. The third one is after social assistance transfers, but before taxes, social insurance, and universal benefits.

By calculating the Gini coefficient and the poverty rate using each of these four income concepts (income before government intervention, income after non-assistance government redistribution, income after social assistance benefits, and income after all government redistribution, i.e. our original disposable household income variable, **dhi**), we gain some insight into the effect of government programmes on inequality and poverty.

*Efficient programming techniques* - This exercise also introduces some programming techniques that allow looping the same code over several variables.

### Activity

As always, begin with the program you developed for the last exercise. Modify it to create three new income variables. The first, **mi**, is the sum of factor income (**hifactor**), private transfers (**hiprivate**) and private pensions (**hi33**). Because we are specifically interested in the role of *government* transfers, we add private transfers and private pensions to our measure of “market income” from labour and capital.

The second, **siti**, adds **mi** together with social insurance transfers (**hpub\_i**) and universal benefits (**hpub\_u**), while subtracting taxes and social contributions paid (**hxitsc**).

The third measure, **sa**, adds together the market income **mi** with social assistance benefits (**hpub\_a**).

The income variable we have been using up to now, disposable household income, adds together the variables contained in **siti** along with social assistance transfers (which are also contained in the variable **hpub\_a**).

Make sure you apply bottom codes and the equivalence scale to the new variables, producing the final variables **emi\_b**, **esiti\_b**, **esa\_b** and **edhi\_b**.

Write a loop to calculate the Gini coefficient and the poverty rate for all four income variables, based on the code from the previous two exercises. Use it to fill out the table below. Make sure you use the *same* poverty line for all four income variables. That is, the poverty line should be defined as 50 percent of the median equivalised disposable household income, and that same poverty definition should be applied to the other three income variables.

	Before taxes and government transfers (mi)	After taxes, social insurance benefits and universal benefits (siti)	After social assistance benefits (sa)	After taxes and all government transfers (dhi)
Gini coefficient				
Poverty rate				

### Question

6.1. Which government instruments have a greater impact on inequality and poverty in Guatemala: taxes/social insurance/universal benefits, or social assistance?

### Guidelines

- Detailed information about income aggregates used in the exercise can be found in [METIS](#): *METIS* → *LIS database* → *select <variables>* → *Results* → *Variable definitions*. In addition to this, it may be useful to consult the current [list of variables](#) on LIS website: *LIS website* → *Our data* → *LIS database* → *list of variables*.
- Note that social assistance, social insurance and universal benefits (**hpub\_a**, **hpub\_i** and **hpub\_u**) together comprise the concept of public transfers. However, the *total sum of these benefits may be lower than the total public transfers reported at the higher-level variable **hpublic***. This is because some amounts that do not clearly belong to one of the transfer categories may be reported directly in the variable **hpublic**. This may also be the case for other higher-level income concepts with an exception of the five blocks of the current income that always add up to the total income (labour income, capital

income, pensions, public social benefits and private transfers). For more information, you can consult the aggregation rules in the [LIS User Guide](#).

- You can easily check whether this is the case by summing up the individual components and the higher-level variable.

## Program

```
gini <- function(x, weight) {
  ox  <- order(x)
  x   <- x[ox]
  weight <- weight[ox]/sum(weight)
  p   <- cumsum(weight)
  nu  <- cumsum(weight*x)
  n   <- length(nu)
  nu  <- nu/nu[n]
  res <- sum(nu[-1]*p[-n])-sum(nu[-n]*p[-1])
return(res)
}

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)
}

setups <- function(data_file) {
  vars <- c('dhi', 'hifactor', 'hi33', 'hpublic', 'hpub_i', 'hpub_u', 'hpub_a', 'hiprivate', 'hxitsc',
'hpopwgt', 'nhhmem', 'grossnet')
  subset <- 'complete.cases(dhi, hifactor, hi33, hpublic, hpub_i, hpub_u, hpub_a, hiprivate,
hxitsc)'
  df <- read.LIS(data_file, labels = FALSE, vars = vars, subset = subset)
  df$dhi <- ifelse(df$dhi < 0, 0, df$dhi)
  df$edhi <- df$dhi / (df$nhhmem^0.5)
  df$mi <- df$hifactor + df$hiprivate+df$hi33
  df$mi <- ifelse(df$mi < 0, 0, df$mi)
  df$emi <- df$mi / (df$nhhmem^0.5)
  df$siti <- df$hifactor + df$hiprivate + df$hi33 + df$hpub_i + df$hpub_u -
df$hxitsc
  df$siti <- ifelse(df$siti < 0, 0, df$siti)
  df$esiti <- df$siti / (df$nhhmem^0.5)
  df$sa <- df$hifactor + df$hiprivate + df$hi33 + df$hpub_a
  df$sa <- ifelse(df$sa < 0, 0, df$sa)
}
```

```

df$esa <- df$sa / (df$nhhmem^0.5)
return(df)
}
df <- setups('gt06h')
maxline <- 0.5
for(var in c('emi', 'esiti', 'esa', 'edhi')) {
  cat(paste("VARIABLE: ", var), sep = '\n')
  cat(paste("Gini Coefficient :", round(gini(df[[var]], df$hpopwgt*df$nhhmem),
digits = 3)), sep = '\n')
  cat(paste("Poverty Rate   :", round(100 * (sum((df[[var]] < maxline *
wNtile(df$edhi, df$hpopwgt * df$nhhmem, 0.5)) * df$hpopwgt * df$nhhmem) /
sum(df$hpopwgt * df$nhhmem))), digits = 2)), sep = '\n')
  cat(" ", sep = '\n')
}
for (var in c('hpublic', 'hpub_i', 'hpub_u', 'hpub_a')) {
cat(paste("Mean: ", var, round(mean(df[[var]]), digits = 2)))
cat(" ", sep = '\n')
}

```

## **Results**

	Before taxes and government transfers (mi)	After taxes, social insurance benefits and universal benefits (siti)	After social assistance benefits (sa)	After taxes and all transfers (dhi)
Gini coefficient	0.518	0.496 *	0.511	0.489
Poverty rate	24.5%	24.0% *	22.9%	22.5%

\* Calculation excludes universal benefits since it is not filled in gt06.

### *Solution*

6.2. Which government instruments have a greater impact on inequality and poverty in Guatemala: taxes/social insurance/universal benefits, or social assistance?

- Taxes, social insurance benefits and universal benefits play a larger role in reducing income inequality in Guatemala in 2006 than social assistance benefits, and the reverse is true for poverty reduction.

### *Comments*

- When comparing incomes before and after taxes and transfers, be careful when interpreting the meaning of the pre-tax and transfer figure. It is tempting to interpret this number as a representation of the income distribution that would exist in the absence of government programs. However, since outcomes in the private sector are conditioned by the presence or absence of government programs, this is not generally a reasonable inference.
- There are several online databases containing detailed information on social security systems around the world that may be useful in your analysis:
  - The Mutual Information System on Social Protection ([MISSOC](#)), which provides up-to-date information on social protection legislation, benefits and conditions for 28 EU Member States, Iceland, Liechtenstein, Norway and Switzerland (and [MISSCEO](#) for countries outside of the EU's MISSOC network);
  - [Social Security Programs Throughout the World](#) – a biannual publication highlighting the principal features of social security programs in more than 170 countries;
  - Other databases and publications from institutional websites, such as the International Labour Organisation ([ILO](#)), the Organisation for Economic

Co-operation and Development ([OECD](#)) and [EUROMOD](#) – the tax-benefit microsimulation model for European Union;

- National social security websites.



## 7. Comparing multiple countries

### Goal

Now that we have written code to compute all of our statistics of interest, it is time to calculate these quantities for multiple countries. Building up on the previous exercise, we will introduce different programming techniques that break the code into logic sub-routines and generalise our program to loop through multiple datasets.

*Adding more countries* – Before adding a new country/year to an analysis, it is important to check that the dataset in question has all information necessary for the analysis you are performing. In this case, one should carefully check that the income sub-components variables used in the previous exercise are filled for all the new datasets to be introduced (and if not, whether the analysis can be slightly modified to take into account a different situation).

*Efficient programming techniques* – In the previous exercise we have introduced some programming techniques that allowed to loop the same code over several variables. In this exercise, we will introduce some other techniques that allow to organise the code in an efficient way and easily loop over both variables and datasets.

### Activity

- (1) Take the code from the previous exercise and modify it so that it loops through five datasets: Guatemala 2006 (**gt06**), United States 2004 (**us04**), Denmark 2004 (**dk04**), Hungary 2005 (**hu05**) and Israel 2005 (**il05**).
- (2) The code that creates the income variables can be placed in a subroutine that is called from the main loop, as can the code that applies the equivalence scale and the bottom code.

Use your results to fill in the following tables:

### **Gini Coefficient**

Dataset	Before taxes and government transfers	After taxes, social insurance benefits and universal benefits	After social assistance benefits	After taxes and all transfers
GT06				
US04				
DK04				
HU05				
IL05				

## Poverty Rate

Dataset	Before taxes and government transfers	After taxes, social insurance benefits and universal benefits	After social assistance benefits	After taxes and all transfers
GT06				
US04				
DK04				
HU05				
IL05				

Keep in mind that even if all cells can technically be constructed, the result may not necessarily be comparable conceptually! Think carefully about whether the dataset you are looking at contains the necessary information to calculate the quantity in each column.

### Question

- 7.1. In what cells does the figure you produced not match the income concept described in the column header?
- 7.2. In which country do government programmes do the most to reduce inequality and poverty, in percentage terms? In which country do they do the least?
- 7.3. In which countries do social assistance benefits do more to reduce poverty than social insurance plus universal benefits and taxes?

## Guidelines

### ➤ Functions

The R language lends itself to functional programming, in which programs are constructed around the inputs and outputs to functions. When a chunk of code needs to be executed repeatedly for different data, it can be useful to put that code in a function, and then have the main program call the function. We have already used functions in the calculation of the Gini coefficient. Since this exercise requires us to perform identical recoding on multiple data sets, we can create a function to do the recoding. The function will look like this:

```
setups <- function(data_file) {  
    <recoding commands from the previous program go here>  
return(df)
```

```
}
```

This function takes a data frame as its argument, performs some recoding on that data frame, and then returns the data frame. So, if you have loaded a dataset into memory, you can then recode it.

➤ Noting irregular cases

The variable **grossnet** reports whether the incomes in a dataset are gross income, before taxes, or whether they only report post-tax values. Within a single dataset, all cases will have the same value for this variable. In datasets containing gross incomes, the **grossnet** variable can take values 111, 110 or 120. In a purely net income dataset, **grossnet** will be 200. Note, however, that a few datasets have **grossnet** codes of 300, 310 or 320 because they contain a mixture of gross and net incomes. More information about whether the LIS datasets report gross or net values can be found on the LIS METIS documentation system at *Our Data* → *METIS* → *LIS database* → *select "ccyy"* → *Dataset information*, while the cross-compared tab allows you to check the variable value taken for the selected dataset(s). METIS is accessed through this [link](#).

➤ Looping through datasets

You can use a loop to iterate over datasets, just as you have used loops elsewhere in your code:

```
for(ccyy in datasets) {  
  df <- setups(ccyy)  
  <other code to produce output>  
}
```

## **Program**

```
gini <- function(x,weight) {
  ox  <- order(x)
  x   <- x[ox]
  weight <- weight[ox]/sum(weight)
  p   <- cumsum(weight)
  nu  <- cumsum(weight*x)
  n   <- length(nu)
  nu  <- nu/nu[n]
  res <- sum(nu[-1]*p[-n])-sum(nu[-n]*p[-1])
return(res)
}

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)
}

setups <- function(data_file) {
  vars <- c('dhi', 'hifactor', 'hi33', 'hpublic', 'hpub_i', 'hpub_u', 'hpub_a', 'hiprivate', 'hxitsc',
'hpopwgt', 'nhhmem', 'grossnet')
  subset <- 'complete.cases(dhi,hifactor,hi33,hpub_i,hpub_u,hpub_a,hiprivate,hxitsc)'
  df <- read.LIS(data_file, labels = FALSE, vars = vars, subset = subset)
  df$dhi <- ifelse(df$dhi < 0, 0, df$dhi)
  df$edhi <- df$dhi / (df$nhhmem^0.5)
  df$mi <- df$hifactor + df$hiprivate+df$hi33
  df$mi <- ifelse(df$mi < 0, 0, df$mi)
  df$emi <- df$mi / (df$nhhmem^0.5)
  df$siti <- df$hifactor + df$hiprivate + df$hi33 + df$hpub_i + df$hpub_u - df$hxitsc
  df$siti <- ifelse(df$siti < 0, 0, df$siti)
  df$esiti <- df$siti / (df$nhhmem^0.5)
  df$sa <- df$hifactor + df$hiprivate + df$hi33 + df$hpub_a
  df$sa <- ifelse(df$sa < 0, 0, df$sa)
  df$esa <- df$sa / (df$nhhmem^0.5)
  return(df)
}
```

```

}
datasets <- c('gt06', 'us04', 'dk04', 'hu05', 'il05')
maxline <- 0.5
for (ccyy in datasets) {
  df <- setups(paste(ccyy,'h',sep="))
  for(var in c('emi', 'esiti', 'esa', 'edhi')) {
    cat(paste("VARIABLE: ", var), sep = '\n')
    cat(paste("Gini Coefficient :", round(gini(df[[var]], df$hpopwgt*df$nhhmem), digits =
3)), sep = '\n')
    cat(paste("Poverty Rate :", round(100 * (sum((df[[var]] < maxline * wNtile(df$edhi,
df$hpopwgt * df$nhhmem, 0.5)) * df$hpopwgt * df$nhhmem) / sum(df$hpopwgt *
df$nhhmem))), digits = 2)), sep = '\n')
    cat(" ", sep = '\n')
  }
print(c(toupper('grossnet')))
print(table(df$grossnet, useNA = 'ifany'))
print(paste(round(prop.table(table(df$grossnet, useNA = 'ifany')) * 100, digits =
2), "%", sep = ""))
cat(paste(" "), sep = '\n')
for (var in c('hpublic', 'hpub_i', 'hpub_u', 'hpub_a')) {
cat(paste("Mean: ", var, round(mean(df[[var]]), digits = 2)))
cat(" ", sep = '\n')
}
}

```

## **Results**

### **Gini Coefficient**

Dataset	Before taxes and government transfers	After taxes, social insurance benefits and universal benefits	After social assistance benefits	After taxes and all transfers
GT06	0.518	0.496 **	0.511	0.489
US04	0.488	0.400	0.468	0.377
DK04	0.422	0.263	0.395	0.230
HU05	0.527 *	0.305	0.510	0.288
IL05	0.502	0.402	0.491	0.377

### **Poverty Rate**

Dataset	Before taxes and government transfers	After taxes, social insurance benefits and universal benefits	After social assistance benefits	After taxes and all transfers
GT06	24.5	24.0 **	22.89	22.5
US04	26.1	21.2	23.0	17.3
DK04	25.3	12.7	22.2	5.6
HU05	42.6 *	9.6	41.3	6.7
IL05	28.6	22.5	27.9	19.4

\*Calculation based on post-tax income

\*\*Calculation excludes universal benefits

### *Solution*

7.1. In what cells does the figure you produced not match the income concept described in the column header?

- Because Hungary 2005 is a net income dataset, the Gini before taxes and transfers cannot be included. While the example program below does produce a result for Hungary, it is not actually comparable to the other countries because it is post-tax. That cell should therefore be left blank.

7.2. In which country do government programmes do the most to reduce inequality and poverty, in percentage terms? In which country do they do the least?

- Government programmes have the largest impact on reducing inequality in Denmark, where they reduce the Gini coefficient by 45 percent, and poverty in Hungary by 84 percent.

- Government programs have the lowest impact on reducing inequality and poverty in Guatemala, where they reduce the Gini coefficient by 6 percent and poverty by 8 percent.

7.3. In which countries do social assistance benefits do more to reduce poverty than social insurance plus universal benefits and taxes?

- In Guatemala, social assistance reduces poverty more than social insurance plus taxes and universal benefits.

#### *Comments*

- The datasets in these exercises were chosen because they allow social assistance to be separated from other kinds of government transfers. In many datasets, unfortunately, this separation is not possible due to the limitations of the original data. In such cases, the total amount of transfer income will be contained in a higher-level variable such as **hpublic**, and lower level variables such as **hpub\_i** and **hpub\_u** will be unfilled. You can consult the **Crossed-compare** function in the METIS documentation tool, to determine which variables are available in each dataset. This is the **main functionality** of METIS. Upon entering the tool ([METIS](#)), it is possible to view the availability of the selected variable(s) in the selected dataset(s) and to **compare information about this selection**. To access this information, *Select datasets* → *Select variables (under Income Aggregates)* → *Results* → *Crossed-Compare*.

## 8. Producing compact and concise output

### Goal

The program we have developed produces results for two indicators (poverty and inequality), four definitions of income, and five countries. This results in a total of 40 values of interest in the resulting log file. We could copy these values into a spreadsheet by hand, but this would be very time-consuming and would increase the likelihood of introducing errors by accidentally copying the wrong number. In this exercise, we will modify the program to create compact output that can be transferred into a spreadsheet with only one cut-and-paste.

The final output will be a set of comma-separated values, in which each country is on a separate row and each indicator is on a separate column.

### Activity

To produce easy-to-use output, we will modify the program as follows:

- Create a matrix to store only the results we need from the program, and label its dimensions appropriately.
- Replace the code that prints out results with code that stores those results in the matrix.
- Print out the matrix at the very end of the program, in Comma Separated Values format.

The result will be tables like the ones shown below. The easiest way to do this is to copy and paste the comma-separated block of results into a spreadsheet

### Guidelines

One advantage of R is that it does not limit the number of separate data objects that can be in memory at one time. This means that we can have a separate matrix to hold our results, while also keeping a LIS data set open for processing. To set up the matrix, put this code before the program's main loop:

```
result      <- matrix(NA,length(datasets),6)  
rownames(result) <- datasets  
colnames(result) <- c('gini_mi','gini_siti','gini_sa','gini_dhi','pov_mi','pov_siti',  
  'pov_sa','pov_dhi')
```

This creates a 5x8 matrix filled with the NA (missing) value, which we will then fill. Inside your main loop, insert the Gini and poverty rate to the appropriate cells.



```
result[match(ccyy,datasets), match(var, c('emi','esiti', 'esa','edhi'))] <- "your-gini"  
result[match(ccyy,datasets), 3 + match(var, c('emi','esiti', 'esa','edhi'))] <- "your-  
pov"
```

By using the **match** command, we can ensure that each number is inserted into the correct column and row. Then at the very end of the program, *after* the main loop, we insert:

```
write.csv(result)
```

This will create output like that shown [below](#).

## **Program**

```
gini <- function(x,weight) {
  ox  <- order(x)
  x   <- x[ox]
  weight <- weight[ox]/sum(weight)
  p   <- cumsum(weight)
  nu  <- cumsum(weight*x)
  n   <- length(nu)
  nu  <- nu/nu[n]
  res <- sum(nu[-1]*p[-n])-sum(nu[-n]*p[-1])
return(res)
}

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)
}

setups <- function(data_file) {
  vars <- c('dhi', 'hifactor', 'hi33', 'hpublic', 'hpub_i', 'hpub_u', 'hpub_a', 'hiprivate', 'hxitsc',
'hpopwgt', 'nhhmem', 'grossnet')
  subset <- 'complete.cases(dhi,hifactor,hi33,hpub_i,hpub_u,hpub_a,hiprivate,hxitsc)'
  df <- read.LIS(data_file, labels = FALSE, vars = vars, subset = subset)
  df$dhi <- ifelse(df$dhi < 0, 0, df$dhi)
  df$edhi <- df$dhi / (df$nhhmem^0.5)
  df$mi <- df$hifactor + df$hiprivate+df$hi33
  df$mi <- ifelse(df$mi < 0, 0, df$mi)
  df$emi <- df$mi / (df$nhhmem^0.5)
  df$siti <- df$hifactor + df$hiprivate + df$hi33 + df$hpub_i + df$hpub_u - df$hxitsc
  df$siti <- ifelse(df$siti < 0, 0, df$siti)
  df$esiti <- df$siti / (df$nhhmem^0.5)
  df$sa <- df$hifactor + df$hiprivate + df$hi33 + df$hpub_a
  df$sa <- ifelse(df$sa < 0, 0, df$sa)
  df$esa <- df$sa / (df$nhhmem^0.5)
  return(df)
}
```

```

}

datasets <- c('gt06', 'us04', 'dk04', 'hu05', 'il05')
maxline <- 0.5
result      <- matrix(NA,length(datasets),8)
rownames(result) <- datasets
colnames(result) <- colnames(result) <-
c('gini_mi','gini_siti','gini_sa','gini_dhi','pov_mi','pov_siti','pov_sa','pov_dhi')

for (ccyy in datasets) {
  df <- setups(paste(ccyy,'h',sep=""))
  # For net income datasets, blank out market income measures, since these cannot
  # be calculated correctly
  for(var in c( 'esiti', 'esa', 'edhi')) {
if (df$grossnet ==200) {
  result[match(ccyy,datasets), match(var, c('emi', 'esiti', 'esa', 'edhi'))] <-
round(gini(df[[var]], df$hpopwgt*df$nhhmem), digits = 3)
  result[match(ccyy,datasets), 4+ match(var, c('emi', 'esiti', 'esa', 'edhi'))] <-
round(100 * (sum((df[[var]] < maxline * wNtile(df$edhi, df$hpopwgt *
df$nhhmem, 0.5)) * df$hpopwgt * df$nhhmem) / sum(df$hpopwgt * df$nhhmem)),
digits = 2)
} else {
  for(var in c('emi', 'esiti', 'esa', 'edhi')) {
  result[match(ccyy,datasets), match(var, c('emi', 'esiti', 'esa', 'edhi'))] <-
round(gini(df[[var]], df$hpopwgt*df$nhhmem), digits = 3)
  result[match(ccyy,datasets), 4 + match(var, c('emi', 'esiti', 'esa', 'edhi'))] <-
round(100 * (sum((df[[var]] < maxline * wNtile(df$edhi, df$hpopwgt *
df$nhhmem, 0.5)) * df$hpopwgt * df$nhhmem) / sum(df$hpopwgt * df$nhhmem)),
digits = 2)
}
}
}
}
print(write.csv(result))

```

## **Results**

If you have modified the program correctly, you should see a block of text like this at the end of your log file:

```
"","gini_mi","gini_siti","gini_sa","gini_dhi","pov_mi","pov_siti","pov_sa","pov_dhi"  
"gt06",0.518,0.496,0.511,0.489,24.48,24.03,22.89,22.49  
"us04",0.488,0.4,0.468,0.377,26.12,21.18,23.01,17.32  
"dk04",0.422,0.263,0.395,0.23,25.35,12.7,22.16,5.64  
"hu05",NA,0.305,0.51,0.288,NA,9.62,41.3,6.72  
"il05",0.502,0.402,0.491,0.377,28.6,22.53,27.95,19.37
```

If you copy and paste this into a spreadsheet, most spreadsheet programs should recognize this as a set of comma-separated values and parse it automatically. You can also copy the text into a text file, save it with the **.CSV** extension, and open it with your spreadsheet program.

For example, you can copy and paste the values into an Excel file, use the *Text To Columns* function under the *Data* tab and choose *comma* as a delimiter.

## 9. Producing graphs with LISSY

### Goal

This exercise introduces LISSY's graphing feature that allows users to generate, display and export graphs based on LIS data. Your goal is to visually compare income inequality across different countries, using our variable for disposable household income (**dhi**).

### Activity

In this exercise you will generate a Lorenz Curve on disposable household income (dhi) for Guatemala (gt06), the United States (us04), and Denmark (dk04) in order to compare the degree of income inequality among these countries.

- (1) Prepare the R session, loading the required packages and functions
- (2) Import the files for Guatemala (**gt06**), the United States (**us04**), and Denmark (**dk04**)
- (3) Prepare the data for the analysis, making sure you drop observations with missing data and consider the winsorised (or bottom- and top-coded) version of the equivalised disposable income
- (4) Generate Lorenz Curves Guatemala, United States, and Denmark using the winsorised version equivalised income

### *Question*

- 9.1. Compare the degree of income inequality among Guatemala, the United States and Denmark using Lorenz Curve.

### Guidelines

- Preparing the session

The first thing that you need to do is to start the R session by loading the required packages and functions. For this exercise you should load four R packages, enabling the system to read all the functions stored within the **lissyrtools** directory.

**library(dplyr)**

**library(magrittr)**

**library(purrr)**

**library(ggplot2)**

Additional documentation on these functions, such as arguments and examples

of use can be found [here](#).

### ➤ Read files

The `read_lissy_files()` function will let you import one or multiple files, specified as a character vector in the first argument (named 'files'). The function returns a list with the imported files as elements. In this exercise, the list is stored as `lissy_datasets` and would contain six elements, one for each imported dataset. The second argument of the function (`full_year_name`) will allow you to obtain the names of the files with four digit years (e.g. `ca2017h` instead of `ca17h`). This can be convenient when performing certain actions, as four digit years can be more easily sorted.

```
lissy_datasets <- read_lissy_files(files = c("gt06h", "us04h", "dk04h"),  
  full_year_name = TRUE)
```

### ➤ Data management

After importing the datasets, you will prepare them for the analysis. For this exercise, you can use an alternative approach based on `lissytools` that offers a set of `transform_` functions such as the ones below:

```
lissy_datasets %<>%  
  transform_negative_values_to_zero(variable = "dhi") %>%  
  transform_equivalise(variable = "dhi") %>%  
  transform_top_code_with_iqr(variable = "dhi", times = 3) %>%  
  transform_weight_by_hh_size(variable = "dhi")
```

The R package `magrittr`, `%>%` and `%<>%` are pipes that allow to pass the result of one function as the first argument of the next one. Additionally, `%<>%` stores the result back to the left-hand-side object passed. Functions in `lissytools` are compatible with these pipes as they make the code much easier to read.

The `transform_` functions above, process the files previously stored as `lissy_datasets` with the following:

- `transform_negative_values_to_zero (variable = "dhi")` recodes all negative values to zero in the selected variable 'dhi' (disposable household income).
- `transform_equivalise (variable = "dhi")` adjusts the selected variable by square root of the number of household members.
- `transform_top_code_with_iqr (variable = "dhi", times = 3)` applies an upper limit to the variable. This corresponds to 3 times the Interquartile Range of the variable transformed using the natural logarithm.
- `transform_weight_by_hh_size (variable = "dhi")` multiplies the weight by

the number of people in the household.

Notice that, as we used the `%<>%` pipe, all transformations are automatically stored back to `lissy_datasets`.

➤ Generating and exporting graphs

You can plot the Lorenz curve for a variable of all imported datasets with the `plot_lorenz_curve()` function. Please make sure that missing values are removed. The function also has an optional `plot_theme` argument, which currently supports only a limited number of options.

```
lissy_datasets %>%
```

```
  plot_lorenz_curve(variable = "dhi", na.rm = TRUE, plot_theme = "lis")
```

## **Program**

```
# Prepare session
library(readr)
library(dplyr)
library(magrittr)
library(purrr)
library(ggplot2)

all_lissyrttools_scripts <- fs::dir_ls("/media/user/lissyrttools/")
invisible(purrr::map(all_lissyrttools_scripts, ~ source(.x)))

# Script -----

# Read files
lissy_datasets <- read_lissy_files(files = c("gt06h", "us04h", "dk04h"), full_year_name =
TRUE)

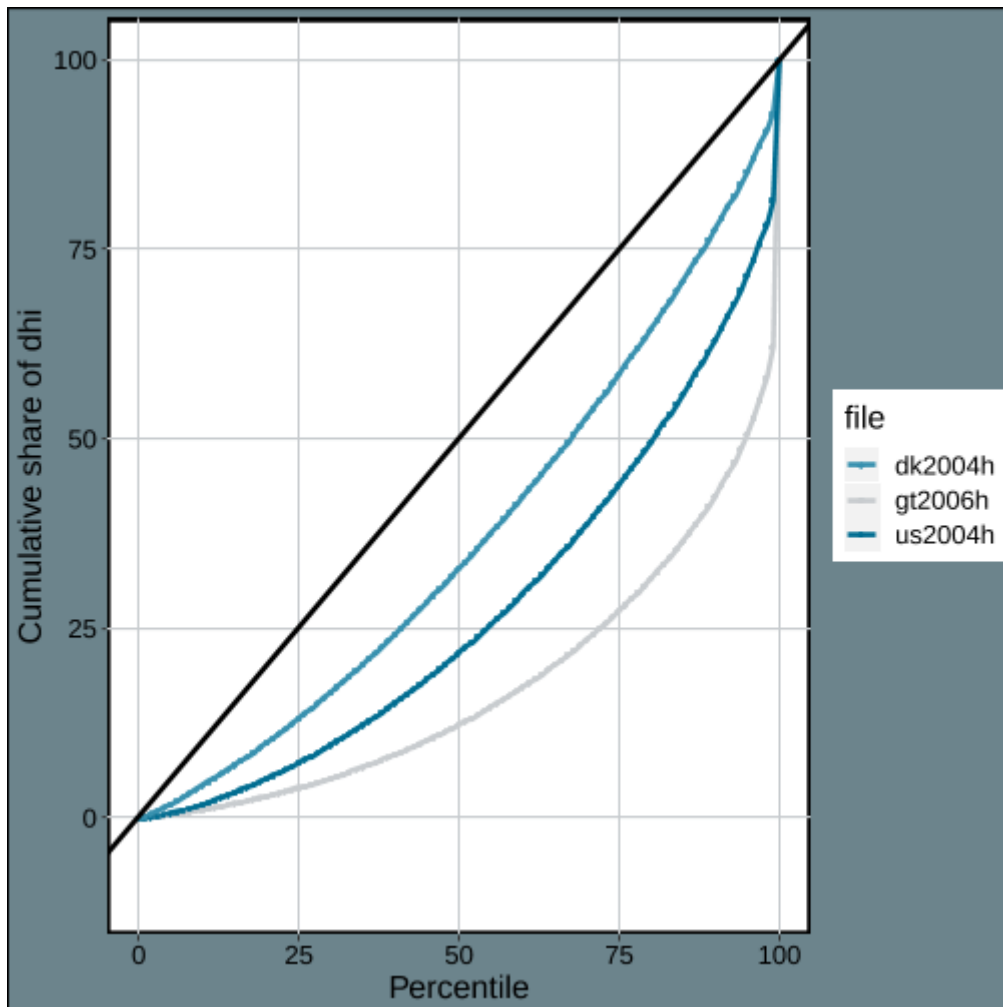
# Data management
lissy_datasets %<>%
  transform_negative_values_to_zero(variable = "dhi") %>%
  transform_equalise(variable = "dhi") %>%
  transform_top_code_with_iqr(variable = "dhi", times = 3) %>%
  transform_weight_by_hh_size(variable = "dhi")

# Plot Lorenz Curve
lissy_datasets %>%
  plot_lorenz_curve(variable = "dhi", na.rm = TRUE, plot_theme = "lis")
```



## Results

**Graph – Lorenz Curve comparison between Guatemala (2006), United States (2004) and Denmark (2004)**



### *Solution*

9.1. Compare the degree of income inequality among Guatemala, United States, and Denmark using Lorenz Curve.

- The Lorenz curve is a classical tool for describing income distributions that allows making pair-wise rankings, in the case in which the curves do not cross.
- Among the countries considered, the United States is the one with highest level of inequality, while Denmark is the one with the lowest.
- It is important to remind that intersecting Lorenz curves prevents conclusions on which income distribution has more inequality.