

# Using R to extract data from the World Bank's World Development Indicators

BY DOMINIQUE VAN DER MENSBRUGGHE<sup>a</sup>

*This article describes three applications which use a new R-based interface that facilitates extraction of data from the World Bank's World Development Indicators (WDI) database. WDI contains over 1300 time series of socio-economic indicators and provides important inputs for the Global Trade Analysis Project (GTAP) Data Base and, as well, can assist in the development of socio-economic scenarios. The first application is a stand-alone R application that extracts per capita historical GDP for all countries and summarizes growth episodes across regions in a box plot. The two other applications highlight the increasing use of interoperability across software packages to take advantage of the specific features of each. The first couples the R-based WDI extraction with a 'Beamer'-type L<sup>A</sup>T<sub>E</sub>X presentation. This is a potentially powerful combination for a quick (and programmatic) turnaround from data to presentation. With little effort, the application can be extended to generate a large number of slides. The second example of interoperability is centered on a General Algebraic Modeling System (GAMS) program that calls R as a sub-process to extract the data for further processing in GAMS. The R-script developed for this last application could be readily coupled with other software packages.*

JEL codes: C10, C82, C88

Keywords: Economic database; Software; R; Latex; GAMS

## 1. Introduction

With a strong push from then-President Robert Zoellick, the World Bank made publicly and freely available its major database known as the World Development Indicators (WDI). WDI contains a large host of socio-economic indicators, from the widely used such as population and GDP, to the more esoteric such as the percent of households that consume iodized salt. At last count WDI contained over 1300 indicators covering most of the countries of the world, and with the earliest indicators starting in 1960. Providing easy access to these indicators significantly cuts the cost of undertaking detailed quantitative analysis. Assessing historical trends—with either time series analysis or structural econometrics—is a *sina qua non* for

---

<sup>a</sup> The Center for Global Trade Analysis, Purdue University, 403 West State Street, West Lafayette, IN 47906 (e-mail: vandermd@purdue.edu).

shaping long-term scenarios of growth, development and sustainability, such as frequently done in the context of climate change analysis, for example.

There are two methods for accessing WDI. The first is to directly use the World Bank's web-based graphical user interface (GUI).<sup>1</sup> The user points and clicks with a mouse to select countries, indicators and years. The selected data can then be downloaded in alternative formats, for example Excel. This access method is fine for the occasional use of WDI but quickly becomes tedious for large selections and/or when access is routine. The second method uses a so-called Application Programming Interface (API) that can be embedded in computer code to programmatically extract data from WDI. The API requires inputs—selected countries, indicators and years—and returns the desired data 'cube'. The API has been integrated into an R package<sup>2</sup> that simplifies the extraction process and allows for the downloaded data to be directly treated in R or to be saved as a datafile for use with another programming environment.

Each new version of the Global Trade Analysis Project (GTAP)<sup>3</sup> database extensively uses data in the WDI database for shaping the macroeconomic framework—this includes indicators such as GDP and its components (private consumption, government and investment expenditures, etc.), population, the stock of capital and additional accounts required for satellite data sets and/or for historical purposes. For example, WDI includes the share of agriculture in GDP, the rate of urbanization, the share of the working age population in total population, energy use and CO<sub>2</sub> emissions, further decomposition of the balance of payments such as remittances, transfers and aid to governments, and poverty and distributional statistics.

This article provides a broad description of the type of data series available in WDI and how to use the new R-based extraction package to extract the data. A first application extracts real per capita growth for all countries since 1960 and plots long-run growth episodes in a box plot where countries are classified using the World Bank's regional definitions. This application illustrates the basic functionality of the extraction routine and how to combine the extraction with statistical analysis in R. It also shows how to use R's plotting features to save figures that can

---

<sup>1</sup> The GUI-interface to WDI is available at the following link: <http://databank.worldbank.org/data/reports.aspx?source=world-development-indicators>.

<sup>2</sup> R is open source software designed specifically for statistical analysis but with significant additional functionality for data and graphical analysis. More information, including how to download the software, is available at <https://www.r-project.org/>.

<sup>3</sup> The GTAP Data Base is a database of global economic transactions that integrates and harmonizes national input-output (IO) tables with a consistent database of bilateral trade flows. The latest release (Version 9) includes IO tables for 120 countries (with all remaining countries aggregated in 20 geographical regions) and economic activity is divided into 57 sectors. The database is described in greater detail in Aguiar et al. in this issue. Additional information on GTAP is available at <https://www.gtap.agecon.purdue.edu/>.

be imported into documents. To highlight the increasingly wide use of interoperability across different types of software packages, two additional applications are developed. The first application links R with L<sup>A</sup>T<sub>E</sub>X's 'Beamer' document class—a PDF slide presentation.<sup>4</sup> It illustrates how to quickly generate a presentation from the extracted WDI data, bypassing the at-times tedious steps of using Excel and/or Powerpoint packages. The second illustration embeds the R extraction package in a General Algebraic Modeling Systems (GAMS)<sup>5</sup> program that lets the user select the countries, indicators and range of years to be extracted. The selection is then passed on to R and the data returned to GAMS for additional processing.

## 2. World Development Indicators

The World Development Indicators database is one of the most significant international databases and is developed and maintained at the World Bank. The indicators stored in the database are largely socio-economic indicators at a national level. The most widely used are probably GDP and its components, and population. The indicators are typically annual time series many of which go back to 1960. A large number of the indicators are developed directly at the World Bank (at times in coordination with other international agencies such as the International Monetary Fund (IMF) and the United Nations (UN) Population Division), but a significant number are integrated from external databases. WDI includes, for example, energy data from the International Energy Agency (IEA), labor data from the International Labor Organization (ILO) and agricultural data from the United Nations Food and Agriculture Organization (FAO). Depending on specific needs, WDI can be used for one-stop shopping rather than having to refer to the more specialized databases. Moreover, with open access to WDI complemented with its API, data extraction can be done at extremely low cost.

The WDI data can be used for a variety of purposes—illustrating historical trends, parameter estimation, calibration of trends, assessing structural changes over time—for example population cohorts, the share of agriculture in GDP and employment, and a number of indicators related the so-called Millennium Development Goals (MDGs).<sup>6</sup>

---

<sup>4</sup> L<sup>A</sup>T<sub>E</sub>X is open source software designed for the preparation of documents. Traditionally it has been used to convert L<sup>A</sup>T<sub>E</sub>X input files into PDF files such as papers, reports, articles and books. More recent extensions include the preparation of slide presentations and web-based documents. L<sup>A</sup>T<sub>E</sub>X is freely available for a number of computing platforms and more information is available at <https://latex-project.org/>.

<sup>5</sup> GAMS is widely used in economics and other disciplines for mathematical programming. It provides an intuitive language for the development of mathematical models for the purpose of optimization and simulation. More information on GAMS, including extensive documentation, is available at [www.gams.com](http://www.gams.com). GAMS is freely available, though its use is limited without acquiring a license.

<sup>6</sup> More recently called the Sustainable Development Goals (SDGs).

Extraction from WDI requires selecting countries and/or regions, indicators, and the years to extract. The GUI provides the easiest access as the extraction is based on point and click and the user does not need to know the underlying arcane nomenclature for the countries/regions and the indicators. However, the API depends on 'codes' and these must be provided with exactitude to extract the desired data. Various systems have been devised for country codes. Perhaps the most widely used system is the ISO-3 coding.<sup>7</sup> For example, the ISO-3 code for Switzerland is CHE.<sup>8</sup> The World Bank country codes hone closely to the ISO-3 codes, but it does not necessarily update codes when they change.<sup>9</sup> The R WDI extraction package is based on a different set of codes known as ISO-2 that rely on 2-letter country codes.<sup>10</sup>

The indicator codes<sup>11</sup> have been developed by the World Bank and are typically short descriptions of the underlying indicator—divided into as many as five fields. The first field reflects the topic. For example, 'SP' refers to social indicators and population. The indicator for total population is SP.POP.TOTL. The GDP indicators are in the broad topic of national income (NY). The second field is the type of national income indicator—for example GDP or GNP. The third field will indicate the type of income indicator, for example per capita (PCAP) or total at market price (MKTP). The fourth indicator will most often indicate the pricing regime—current local currency, dollar, purchasing power parity (PPP) dollars, or at constant prices at some base year. For example, GDP at constant dollars at market exchange rates is given by NY.GDP.MKTP.KD.<sup>12</sup>

### 3. A practical illustration of the WDI extraction routine in R

This section describes an R script that plots long-run growth episodes since 1960 across all countries organized into the World Bank's regional classifications.

Figure 1 represents a box plot of all 10-year country growth episodes from the WDI database from 1960 through 2014 with the growth episodes classified according to the World Bank's classification: East Asia & Pacific (EAP), South Asia (SAS), Europe & Central Asia (ECA), Middle East & North Africa (MNA), Sub-Saharan Africa (SSA) and Latin America & Caribbean (LAC). The six geographic regions only include developing countries using the World Bank's definition. All of the

---

<sup>7</sup> See [https://en.wikipedia.org/wiki/ISO\\_3166-1\\_alpha-3](https://en.wikipedia.org/wiki/ISO_3166-1_alpha-3).

<sup>8</sup> Coming from the Latin *Confoederatio Helvetica* or *Confédération Helvétique* in French.

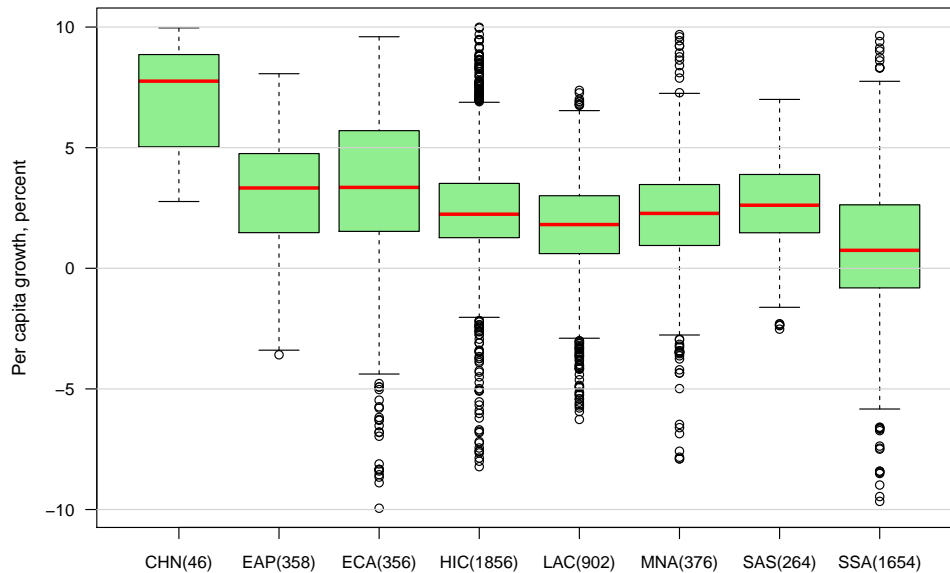
<sup>9</sup> For example, the World Bank continues to use the old ZAR code for the Democratic Republic of the Congo instead of the ISO-3 code of COD. Similarly for Romania, the ISO-3 code is ROU, and not the World Bank's ROM.

<sup>10</sup> See [https://en.wikipedia.org/wiki/ISO\\_3166-1\\_alpha-2](https://en.wikipedia.org/wiki/ISO_3166-1_alpha-2). Note that these codes are widely used for Internet country suffixes.

<sup>11</sup> Also known as Catalog of Economic Time Series or CETS codes.

<sup>12</sup> A mostly complete list of CETS and country codes is available from the author.

developed economies are included in HIC, i.e. the high-income country classification. China has been assigned to its own region (CHN) and excluded from the EAP region. The maximum number of 10-year growth episodes for any single country over the period 1960-2014 is 46.<sup>13</sup> The red line represents the median for the given region. China has the highest, with a median of over 7.5 percent. The HIC region has a median of around 2.5 percent but with a significant number of outliers on both the low- and high-sides.<sup>14</sup> Note that the HIC region now includes a number of European transition economies that experienced a sharp decline after 1990. Not surprisingly, Sub-Saharan Africa shows the lowest median, at around 1 percent, with an inter-quartile range that includes the origin.



**Figure 1.** Growth periods across regions (count in parenthesis).

Source: Author calculations.

These numbers suggest some possibilities for projecting forward. The inter-quartile ranges could be used to bound low- and high-growth scenarios—perhaps attached to specific storylines. An alternative approach is to use the underlying probability distribution to bootstrap growth episodes for the future.

The box plot is created in R using a script that uses the WDI package to directly

<sup>13</sup> The historic time span, the length of a growth episode and the outlier cutoff are all user-defined options. Herein the time span has been set to 1960-2014, the length of a growth episode is 10 years and the outlier cutoff has been set to 10 percent.

<sup>14</sup> Outliers are represented by hollow circles in the chart.

**Table 1.** User options.

Name	Description
wDir	Path for the working directory. This directory will contain both the input file and the box plot figure. To use the default directory, use './'.
fName	Name of the input file that contains the country codes to extract and their associated region.
bpName	Name of the output file containing the box plot figure.
fType	The file type for the box plot. The choices are PDF, JPG and WMF (Windows Metafile format).
begYear	First year of extraction (the earliest is 1960).
endYear	The final year of extraction.
gap	The length of a growth episode.
threshold	Growth episodes with a growth rate (in absolute terms) greater than the threshold are dropped. This affects the span of the box plot. If the threshold is set to a high number, the number of outliers increases.

Source: Author.

extract the relevant data. The script performs the following actions:

- It loads a short file with all of the country codes and the countries classified according to the World Bank's regional definitions as described above. In total there are 8 regions, with an unbalanced number of countries in each.
- For each country in each region, a growth episode is defined over a fixed period—given by a (user-specified) 10-year span. For each ten-year time span, the 'average' growth is calculated using a log-linear regression of real per capita income as a function of time:  $\ln(y) = \alpha + \beta t$ , where  $y$  is per capita GDP (in constant dollars) and  $t$  is time. The parameter  $\beta$  is the 'average' estimated growth rate and is saved as the value of a given growth episode. For a given time range,  $t1$  through  $tn$ , there will be  $tn - t1 - span + 2$  growth episodes, where  $span$  is the span of a growth episode, for example 10 years. A range of 1960 through 2014 would lead to 46 periods of 10 years—assuming no missing data. The maximum number of growth periods for each region is then  $N(tn - t1 - span + 2)$ , where  $N$  is the number of countries in the region.
- The growth episodes are stored in a matrix with  $R$  columns, where  $R$  is the number of regions and the rows contain the individual growth episodes for all countries in the region for all possible growth episodes.
- The growth episodes are plotted in a box plot that is stored in a graphics file that can subsequently be inserted in a document.

The complete R script is provided in Listing B.1.<sup>15</sup> For those familiar with R, most of the code should be relatively self-explanatory. The user's input is handled in the first part—though line 31. Table 1 explains the key options the user can modify.

The second part initializes the script. It loads the WDI library (see Appendix A on how to install the WDI package). It next reads the country codes to be extracted. The codes are stored in a comma separated CSV file. The first line of the file contains a header line. Each subsequent line contains four fields. The first field is a sequence number (that is ignored). The second and third columns contain respectively the ISO-3 and ISO-2 country codes. And the fourth column contains a region acronym. In our example, there are 8 regions as described above. The lines below show an example of the first few lines of what the CSV file could look like:

```
Seq, ISO3, ISO2, WBREG
44, CHN, CN, CHN
82, FJI, FJ, EAP
107, IDN, ID, EAP
123, KHM, KH, EAP
130, LAO, LA, EAP
163, MMR, MM, EAP
166, MNG, MN, EAP
172, MYS, MY, EAP
192, PHL, PH, EAP
194, PNG, PG, EAP
236, THA, TH, EAP
256, VNM, VN, EAP
5, ALB, AL, ECA
```

Line 47 creates a frequency table with the unique region identifiers. R's `table` function is used to build the frequency using column 4 as the key. The table has two columns. The first contains the unique region identifiers and the second the frequency count of countries in the region. Line 50 then determines the number of regions. Line 54 determines the maximum number of countries in any given region. This information is used to calculate the maximum number of possible growth episodes in a region.

Starting in line 64, there is a lengthy loop over all of the regions. Line 67 extracts all of the ISO-2 codes that have been assigned to the region. Lines 71 and 72 perform the critical extraction from the WDI database. The extraction uses the `WDI` function that comes with the WDI R package. The first argument is the vector containing the ISO-2 codes. The second argument is a list of one or more indicators to extract.<sup>16</sup> The final two arguments are the beginning and end years for the extraction.

Once the data is extracted, the growth episodes are calculated for all countries in the region (line 80). The extracted data is ordered by year in ascending order

---

<sup>15</sup> All the code listings described in this article are detailed in full in Appendix B.

<sup>16</sup> The example extracts one indicator at a time, though the WDI function is capable of extracting multiple indicators at a time.

**Table 2.** Summary table for box plot application.

<b>Name</b>	<b>Description</b>
<b>Purpose</b>	The R-based WDI extraction routine is used to extract and plot historical per capita GDP.
<b>Input files</b>	
gdppcplot.R	R-script containing instructions for data extraction and saving box plot chart.
wbreg.csv	User configured input file containing countries to extract and their classification by region. (User can change input file name.)
<b>Output file</b>	
gdppcBoxPlot.pdf	Box plot of per capita income growth episodes. (User can change output file name.)
<b>Instructions</b>	
	Copy the two input files to a user-specified directory. Change the user inputs in the file <code>gdppcplot.R</code> —particularly the directory information. Run the R-script. In R, type <code>source("userDirectory/gdppcPlot.R")</code> replacing <i>userDirectory</i> with the correct path to the directory containing the R script (unless it is in a default directory in which case it can be dropped). The default directory can be set by using the command <code>setwd("userDirectory")</code> .
<b>Software requirements</b>	
R	R statistical software.

Source: Author.



and per capita GDP is calculated (line 91). Growth rates are then calculated for each possible growth episode (line 99). Working vectors  $y$  and  $x$  are filled with the log of per capita GDP and time, respectively. If the sample is full, i.e. there are no missing data, a linear regression is performed to estimate the growth for the episode (line 115). R's `lm` function is used to fit linear models. In this case, the natural logarithm of per capita GDP ( $y$ ) is regressed on time ( $x$ ).<sup>17</sup> The `coef` function is used to extract the second coefficient, i.e. the growth rate, which is converted to a percent growth rate. The coefficient is kept if it is within the user-provided threshold.

At the end of the loop for each region, the vector column of growth rates for the growth episodes is concatenated to the previous (regional) column (line 137), except of course for the first column. The regional acronyms are assigned as the column names. The length of the columns will vary across regions. However, they are padded out by setting the missing cells to NA.

The final part of the script creates the box plot and saves it to a file using the desired format.

#### **4. Interoperability Case 1: WDI extraction and a L<sup>A</sup>T<sub>E</sub>X 'Beamer' presentation**

This section illustrates how to quickly produce a presentation of multiple (line) charts based on an extraction from WDI using R. A single R script is used to select the countries/regions, indicators and time span for using the R-based WDI extraction. The extracted data is then saved as a L<sup>A</sup>T<sub>E</sub>X file that will create a 'Beamer'-type presentation that includes a title slide. In generalized form, this script can be used to generate tens, if not hundreds of slides, with a few short key strokes, and could also be converted to produce other types of charts such as bar charts. L<sup>A</sup>T<sub>E</sub>X has been around for decades and is a very powerful markup language for producing documents, web pages and presentations. It is available for free on many different platforms and has a large user group. This application is relatively self-contained, i.e. it can be used with very little knowledge of L<sup>A</sup>T<sub>E</sub>X, though it requires L<sup>A</sup>T<sub>E</sub>X processing to convert the L<sup>A</sup>T<sub>E</sub>X input into a PDF presentation. Listing B.2 provides the R-script that will produce the L<sup>A</sup>T<sub>E</sub>X based input.<sup>18</sup> Virtually all of the user input is contained in the first part of script (between lines 10 and 90). Table 3 describes the user inputs available for the R-script. The descriptions should be relatively self-explanatory. The presentation will have one slide per extracted indicator (plus the title slide). Each chart will contain one trend line per extracted region, hence it is best to limit the number of countries/regions to extract.<sup>19</sup> Somewhat embedded in the script, the user can also change the theme and theme color of the Beamer

---

<sup>17</sup> By default, the `lm` function includes an intercept term.

<sup>18</sup> L<sup>A</sup>T<sub>E</sub>X files typically have the extension `.tex`.

<sup>19</sup> The script could readily be expanded to chart a much larger set of countries by dividing them in appropriate categories.

presentation.<sup>20</sup>

After the user options, the script initializes a number of variables and loads the WDI library. The preamble for the  $\LaTeX$  file starts at line 115. The  $\LaTeX$  package *pgfplots* will be used to generate the line charts. Line 128 ends the preamble, which produces the title slide. Subsequently, the script loops over all chosen indicators and will produce one slide with a line chart for each indicator. The loop starts in line 132. The extraction from WDI occurs first (line 136). Each chart has a preamble that is contained in lines 140-155. The information used in the preamble is the description of the indicator for the slide title and the beginning and end years for the x-axis. The following section loops over all of the regions for which data has been extracted. Each vector is extracted from the cube (line 170), sorted by year (line 173), and then plotted (lines 180-187) after scaling (line 184). The legend is added at line 187. The slide is terminated with lines 192-196.

---

<sup>20</sup> *CambridgeUS* is the default theme with a default color of *orchid*. Note that the presentation preamble and the charts contain other values that users are free to modify and test.

**Table 3.** User options for creating the WDI-based Beamer presentation.

<b>Name</b>	<b>Description</b>
outwd	Path for the output directory. The $\LaTeX$ ready file will be saved to this directory. To use the default directory, use './'.
fName	Name of the $\LaTeX$ file (with .tex extension).
pTitle	Title for title page of presentation.
shortTitle	Short title of presentation—will be used for footer on slides.
author	Author's name.
nickName	Abbreviation for author's name—will be used for footer on slides.
address	Address—only used on title page. New lines are given by '\\', but must be in pairs because they need to be 'escaped' in R.
institute	Short name for institute, will become part of the footer on slides.
begYear	First year of extraction (the earliest is 1960).
endYear	The final year of extraction.
regions	This contains the regions for which the data will be extracted. It has one row and three columns per region. The first two columns are the ISO-2 and ISO-3 codes respectively. The ISO-2 codes are used for the extraction and the ISO-3 codes for labeling the lines (and legend). The user can change the ISO-3 codes as they are not used for the extraction. The third column is a long-name for the region. It is currently ignored.
indTable	This contains the indicators that will be extracted from WDI. It has one row and four columns per indicator. The first column has a user-specified short name for the indicator and is currently ignored. The second column contains the World Bank CETS code for the indicator to extract. The third column contains a scaling factor for the extracted data. The fourth column contains a long-name for the indicator. This will be used as a slide title.
Colors	A list of colors for the lines on the charts. The list can be as long as desired. A full list of colors is available with the $\LaTeX$ <i>xcolor</i> package.

Source: Author.

**Table 4.** Summary table for Beamer application.

<b>Name</b>	<b>Description</b>
<b>Purpose</b>	The R-based WDI extraction routine is used to extract data series for a number of countries/regions. The data is saved in $\LaTeX$ format to create a Beamer-style presentation of the extracted series as line charts.
<b>Input file</b>	
extractWDIBeamer.R	R-script containing instructions for data extraction and saving data for input into $\LaTeX$ to create a Beamer-style presentation.
<b>Intermediate files</b>	
pptExample.tex	$\LaTeX$ file containing instructions for creating Beamer-style presentation of extracted data. (User can change file name.)
<b>Final file</b>	
pptExample.pdf	PDF presentation of extracted data including title slide. (User can change file name.)
<b>Instructions</b>	
	Copy the file <code>extractWDIBeamer.R</code> to a user-specified directory.
	Change the user inputs in the file <code>extractWDIBeamer.R</code> .
	Run the R-script. In R, type <code>source("userDirectory/extractWDIBeamer.R")</code> replacing <b>userDirectory</b> with the correct path to the directory containing the R script (unless it is in a default directory in which case it can be dropped). The default directory can be set by using the command <code>setwd("userDirectory")</code> .
	The R-script will have created the input file for $\LaTeX$ . Open the $\LaTeX$ input file in $\LaTeX$ and process it (this will be specific to the $\LaTeX$ installation). If there are no errors, $\LaTeX$ will have created the presentation as a PDF file.
<b>Software requirements</b>	
R	R statistical software.
Latex	$\LaTeX$ software for converting $\LaTeX$ markup files into a PDF file.

Source: Author.

## 5. Interoperability Case 2: Combining GAMS with the R-based extraction routine

The second case of interoperability describes how to combine GAMS and R to extract data and read it back into GAMS for additional analysis. Other packages similar to GAMS may have the same tools to link to other packages—for both input and output. The driving software is a GAMS program where the user selects the extraction options—countries/regions, indicators and range of years. The GAMS program saves these options in a file that is subsequently loaded by an R-script that is called from within the GAMS program as a sub-process and that executes the extraction and saves the data. The R sub-process stores the data in a CSV-formatted file. The CSV-formatted file is converted to the GAMS Data eXchange (GDX) format that is read by the GAMS program.<sup>21</sup> The GAMS program continues after reading in the extracted data.

The example GAMS program is provided in listing B.3. The code requires an auxiliary file, called `isocodes.gms`, which contains the set of ISO-2 and ISO-3 codes and their corresponding mapping.

The GAMS code is divided into 5 parts. The first part is a series of global options set by the user. These are described in detail in the preamble of the code. The second part is where the user enters information regarding which countries and indicators to select.<sup>22</sup> The set `c2` contains the list of countries to extract—using ISO-2 codes. It is ignored if the global option `'EXTALL'` is set to `true`. The indicators to extract require three set definitions—the user name of the indicators (that can be the same as the World Bank mnemonics), the World Bank mnemonics (or CETS code) of the indicators and then a mapping between the two sets. These are defined respectively in the sets `'var'`, `'wbind'` and `'mapv'`. The parameter `'scale'` is not needed for the extraction, but is used after the extraction to display the data.

The third part of the code saves the user input in a file that will be used as input to the R script. The format of the output file is a CSV file with each row containing a descriptor and a value.<sup>23</sup> The file name of the input file for the R script is set by the user in the global variable called `'BASENAME'`.

The fourth part of the code uses the `execute` command in GAMS to run R as a sub-task of GAMS using R's *RScript* program. *RScript* works just like R but takes all of its inputs from an R script file and it never opens the graphical user interface of R. It works more or less silently and when finished, the resulting extracted data

---

<sup>21</sup> This additional step is necessary because CSV files are read at compile time. However, the extraction is done at execution time. GDX files are read at execution time and thus this extra step solves the sequencing issue.

<sup>22</sup> The years to extract are selected in the global options section.

<sup>23</sup> Note that the file and the R script are most likely not bullet proof and more error checking would be needed. For example, the R script initializes the number of indicators relative to the number of lines it reads—thus blank lines could affect the code.

will be saved in a CSV file. The next step converts the CSV-formatted file into GDX formatted file.<sup>24</sup> A utility, developed at Wageningen University in the Netherlands, is used to convert the R-extracted CSV file into GAMS' GDX format—it is called `CSV_GDX_Tools`.<sup>25</sup>

The fifth part of the code reads the data from the GDX file into the parameter called `wdiData`. Once the data has been read, it can be used for subsequent processing by GAMS. In the listing file, the data is simply scaled and then displayed.

The rest of this section describes the R-script that is coupled with the GAMS program. It is a generic script in the sense that it not dependent on the calling program—in this case GAMS—, i.e. it could be coupled with any other package that has the ability to write the input file for the R-script and read the output CSV file (with or without the optional header line).<sup>26</sup>

The R-script is divided into two parts. The first part is an R function that performs the routine extraction, does some (slight) checking of the extracted data and saves the data in a CSV file. The function uses two imported functions from the WDI package—the `WDI` extraction function (already described above) and the `countrycode` function that is used to convert the ISO-2 codes to ISO-3 codes.<sup>27</sup>

Listing B.4 contains the R code defining the function 'WDI.Retrieve'. The first executable line of the function performs the extraction using the `WDI` extraction function. The function 'WDI.Retrieve' is written to retrieve only one indicator at a time. The purpose of this is so that the resulting CSV file is essentially 'vectorized' with each extracted indicator appended to the previous one. The 'country' argument in the function call contains the list of countries to extract (using ISO-2) codes or alternatively the text `all` in which cases data for all of the countries is extracted. The 'indicator' argument contains a single indicator code using the World Bank's mnemonics. The 'start' and 'end' arguments represent the first and final years to extract.

The rest of the 'WDI.Retrieve' function essentially performs a series of adjustments to the extracted cube. The first adjustment appends the ISO-3 codes to the

---

<sup>24</sup> GAMS users could bypass directly the CSV formatted-file and instead use the `gdxrrw` package that has been developed to read and write GDX files in R. We have kept this version as the CSV-formatted file can be useful in many other contexts such as Stata, Excel, etc. The `gdxrrw` package for R can be downloaded from [https://support.gams.com/gdxrrw:interfacing\\_gams\\_and\\_r](https://support.gams.com/gdxrrw:interfacing_gams_and_r).

<sup>25</sup> Note that this data conversion utility also requires the file 'forbiddenwords.txt' that should be downloaded with the utility. Both files should be stored in a folder that is accessible (through the system path). Since it is meant to be coupled with GAMS, the GAMS folder would be a logical place to put the utility. The software is available for download at <http://www3.lei.wur.nl/gamstools/index.htm>.

<sup>26</sup> The author has coupled this script with a VBA macro in Excel that facilitates the choice of indicators and reads the subsequent CSV file into an Excel pivot table.

<sup>27</sup> The `countrycode` function comes in its own package and must be installed with the WDI package. Additional details are available in Section A.

**Table 5.** Summary table for GAMS application.

<b>Name</b>	<b>Description</b>
<b>Purpose</b>	The GAMS based program prepares a file with WDI selection options that are passed to R as a sub-process. The R-extracted data is saved as a CSV file that is converted to the GDX format. The latter is read into GAMS for further processing.
<b>Input files</b>	
gdpr.gms	The GAMS program that calls for the R extraction and then processes the extracted data.
isocodes.gms	A GAMS file that contains set definitions for the ISO-2 and ISO-3 country codes and their mapping. Also contains a mapping of the ISO-3 codes to the GTAP regions.
extractWDI.R	R script that reads the options for the extraction, extracts each indicator and saves the extracted data to a CSV file.
<b>Intermediate files</b>	
gdpr.opt	File containing the WDI selections for the R-based extraction.
gdpr.csv	File containing the data extracted by R.
gdpr.gref	File created by the CSV to GDX conversion tool. It is ignored.
<b>Final file</b>	
gdpr.gdx	The extracted data in GDX format
<b>Instructions</b>	
	Copy the files <code>gdpr.gms</code> , <code>isocodes.gms</code> and <code>extractWDI.R</code> to a user-specified directory.
	Change the user inputs in the file <code>gdpr.gms</code> .
	Run the GAMS program <code>gdpr.gms</code> . If using the GAMS Integrated Development Environment (IDE), open the file <code>gdpr.gms</code> and run it.
	If the program runs without error, the extracted data will be printed in the GAMS listing file and is saved in a GDX file.
<b>Software requirements</b>	
GAMS	GAMS mathematical programming software.
R	R statistical software.
CSV_GDX_Tools	Data conversion package that converts CSV files to GDX format (and vice versa).

Source: Author.

data cube using the `countrycode` translation function. Not all ISO-2 codes have an ISO-3 code. The second adjustment allows for the replacement of ISO-3 codes—existing or not. It is a two-step procedure. In the first step, the column indices of the ISO codes are located using the 'grep' (or search function). The second step replaces the relevant ISO-3 codes for the selected ISO-2 codes. In the listing example two ISO-3 codes are replaced (or created)—for Jersey and Guernsey Islands and for Kosovo.<sup>28</sup> The third adjustment is to delete all rows with an unknown ISO-3 code. The fourth adjustment appends the user-specified name of the indicator to each row of the cube. The user-specified name is typically a shorter name of the indicator compared to the World Bank mnemonic (or CETS code). For example `pop` might be used for total population rather than `SP.POP.TOTL`. The final adjustment re-arranges selected columns of the cubes. The ISO-2 column is dropped and the country name column is dropped if the `ifName` flag is set to *false*. Typically the country name can be retrieved from a lookup table given the ISO-3 code and its output in the data cube is redundant and space consuming. The final step is to save the data. The output is a CSV-formatted file. Each row contains the name of the indicator, the ISO-3 country code, the year and the value of the indicator. The country name is an optional field (see above).

The 'WDI.Retrieve' function is relatively generic and can be used multiple times in a single R script. In the scripts below, it is called for each individual indicator.<sup>29</sup>

The R code in listing B.5 is used to read the input file prepared by GAMS and to extract the selected data. It will call the 'wdi.Retrieve' function described above.

The first set of executable lines load the *WDI* and *countrycode* packages from the library. The next section of the code opens the GAMS-prepared file with the options and reads and parses each option. These include the name of the output file, the first and final years, a true/false flag for the header, a true/false flag to output or not the country name in the CSV file and a comma delimited string that contains the countries to be extracted (or else 'all' if all countries are to be extracted). The parsing functions use a few tricks in R. The 'grep' function is used to find the name of the relevant option in the file, i.e. the options can be in any order in the file (except for the indicators that must be the final lines in the file). The 'as.numeric' and 'as.logical' functions are used to convert the options into numerical or boolean values. The 'toupper' function is used to ensure that the ISO-2 country codes are in upper case (a requirement of the *WDI* routine), but it checks first to see if 'all' countries have been selected. The final part of the parsing reads the indicators that are the last lines in the file. The indicators are loaded into a matrix called

---

<sup>28</sup> This adjustment is typically only relevant when `all` countries are selected. There is also no indication that the World Bank will continue to maintain these 2-letter ISO codes or add new ones not yet officially recognized.

<sup>29</sup> There might be some cost in speed as the *WDI* function will be called for each separate indicator rather than doing one call with all of the indicators. However, it is easy to create a CSV cube using this methodology.



'indTable'.

The next section of the R script initializes the CSV file—only in the case a header row has been requested. The final part of the R script loops over each of the indicators and uses the 'wdi.Retrieve' function to retrieve each indicator, one at a time, and save them in the CSV file.

### **Acknowledgements**

The author would like to thank Angel Aguiar and Erwin Corong for beta-testing the software. Software advice was provided by Steven Dirkse from the GAMS Development Corporation and Wietse Dol from Wageningen University. Very useful comments were also provided by the editors and two anonymous reviewers.

## Appendix A. Installing the R packages and their features

The R scripts rely on two packages developed by Vincent Arel-Bundock. The first is called *WDI* and it contains the actual WDI extraction routine. The second is called *countrycode*. It contains lists of country codes developed by different agencies—such as the World Bank, United Nations and others. The *countrycode* function translates one coding system into another. The *WDI* function returns the data cube using ISO-2 character codes.<sup>30</sup> These are converted to ISO-3 character codes before saving the data.

Installation of the packages is only required once.<sup>31</sup> To install the packages from within R, enter the following commands:

```
install.packages("WDI")
install.packages("countrycode")
```

A full description of the functionality of both packages can be downloaded with the packages.<sup>32</sup>

This section highlights some of the other features of the *WDI* package. The *WDI* function has an optional parameter, 'extra', that is a Boolean parameter taking TRUE and FALSE values. If the parameter is set to TRUE, the data cube will contain additional country-specific information such as the ISO-3 code, the name of the World Bank region and the income status of the country. The *WDIcache* function will return the most recent list of indicators available in WDI. The *WDIsearch* function can be used to search for indicators containing a specific string, for example 'gdp'. The returned cache contains a wealth of other information such as capitals and latitude and longitude. The following R commands show how to save the information in the returned cache in CSV files that can readily be loaded into an Excel workbook. The cache is returned as two (text) matrices. The first contains the complete list of indicators with codes, description and source database. The second contains the full list of countries and regions with sundry additional information.

```
Cache <- WDIcache() %
write.table(Cache[1], "indTab.csv", row.names=TRUE, sep=", ")
write.table(Cache[2], "regionTab.csv", row.names=TRUE, sep=", ")
```

---

<sup>30</sup> The function can optionally also return the ISO-3 codes, so the use of the *countrycode* function is not strictly needed, but it is maintained to illustrate its functionality and users can modify the R scripts below to convert to different coding conventions.

<sup>31</sup> The *WDI* package also requires the *RJSONIO*, which is normally automatically loaded when the *WDI* package is initialized.

<sup>32</sup> Available at <https://cran.r-project.org/web/packages/WDI/WDI.pdf> and <https://cran.r-project.org/web/packages/countrycode/countrycode.pdf>.

## Appendix B. The program listings

### Listing B.1. Plotting historical growth periods (*gdppcPlot.R*).

```
1 # User data
3 # Set the input/output file names
5 # Windows
6 wDir <- "v:/data/wdi/"
8 # Mac
9 # wDir <- "~/Documents/Dominique/data/wdi/"
11 # Name of input file with country and region codes
12 fName <- "wbreg.csv"
14 # Output file name
15 bpName <- "gdppcBoxPlot"
17 # Type of output for box plot (valid options are PDF, WMF, JPG)
19 fType <- "pdf"
21 # First year
22 begYear <- 1960
24 # Final year
25 endYear <- 2014
27 # Growth interval (in years)
28 gap <- 10
30 # Threshold for dropping growth rates (in percent)
31 threshold <- 10.1
33 # END OF USER SECTION
35 # Initialize files
36 fName <- paste(wDir, fName, sep="")
37 bpName <- paste(wDir, bpName, sep="")
38 fType <- toupper(fType)
40 # Load the WDI library
41 library(WDI)
43 # Read the countries/regions to analyze
44 wb <- read.csv(fName, sep=",", header=TRUE, na.strings="")
46 # Get the codes for the WB regions
47 wbReg <- as.data.frame(lapply(wb, function(x) (table(x))[4]))
49 # Get the number of WB regions
50 nReg <- nrow(wbReg)
52 # Get the greatest number of countries in any given region
53 # And calculate the maximum number of growth episodes for any region
54 maxC <- max(wbReg$WBREG.Freq)
55 maxLen <- maxC*(endYear-begYear-gap+2)
```

```

57 # Initialize the data vectors for the growth rate estimations

59 y <- vector(mode="double", length=gap)
60 x <- vector(mode="double", length=gap)

62 # Loop over all WB regions

64 for(r in 1:nReg){

66 # Get the country ISO-2 codes for this region and get the number of countries
67 iso <- wb[wb$WBREG==wbReg[r,1],3]
68 nc <- length(iso)

70 # Extract the data from WDI
71 pop <- WDI(iso, "sp.pop.totl", start=begYear, endYear)
72 gdp <- WDI(iso, "ny.gdp.mktp.kd", start=begYear, endYear)

74 # Initialize the growth episode vector
75 gr <- vector(mode = "double", length=maxLen)
76 gr <- NA
77 smpl <- 0

79 # Loop over all countries in this region
80 for(c in 1:nc) {

82 # Extract the data for this country from the regional database
83 popc <- pop[pop$iso2c==iso[c],c("iso2c", "sp.pop.totl", "year")]
84 gdp <- gdp[gdp$iso2c==iso[c],c("iso2c", "ny.gdp.mktp.kd", "year")]

86 # Sort the data by year
87 popc <- popc[order(popc$year),]
88 gdp <- gdp[order(gdp$year),]

90 # Calculate per capita GDP
91 gdppc <- gdp[,2]/pop[,2]

93 # Calculate length of this vector
94 maxYears <- length(gdppc)

96 # For every available growth episode of size gap, calculate the growth rate
97 # using OLS

99 for(t in gap:maxYears) {
100 i <- 0
101 t1 <- t-gap+1
102 t2 <- t

104 # Fill the x and y vectors for OLS
105 for(tt in t1:t2) {
106 if(!is.na(gdppc[tt])) {
107 i <- i+1
108 y[i] = log(gdppc[tt])
109 x[i] = i
110 }
111 }

112 # Calculate the growth rate if we have a full vector
113 if(i == gap){
114 # Do the regression and extract the growth rate
115 grRate <- 100*coef(lm(y ~ x))[2]

```

```
116 #           Only use this growth episode if growth is below the threshold
117 #           in absolute terms
118           if(abs(grRate) <= threshold) {
119               smpl <- smpl+1
120               gr[smpl] <- grRate
121           }
122     }
123 }
124 }
125 # Set the name for this region and append the number of growth episodes
126 regName <- paste(toString(wbReg[r,1]), "(", smpl, ") ", sep="")

128 # To concatenate the vectors, they have to be of identical length
129 # Fill out the vector with NA
130 s1 <- smpl+1
131 for(z in s1:maxLen) {
132     gr[z] = NA
133 }

135 if(r == 1) {
136     result <- cbind(gr)
137 }
138 else {
139     result <- cbind(result, gr)
140 }
141 colnames(result)[r] <- regName
142 }

144 # Create the box plot
145 if (fType == "PDF") {
146     pdf(paste(bpName, ".pdf", sep=""), height=6.5, width=9)
147 } else if (fType == "WMF") {
148     win.metafile(paste(bpName, ".wmf", sep=""), height=5.5, width=8.5)
149 } else if (fType == "JPG") {
150     jpeg(paste(bpName, ".jpg", sep=""), height=480, width=580)
151 } else {
152     stop(paste("Wrong file type <", ftype, ">", sep=""))
153 }

155 # Create the box plot

157 par(cex.axis=0.75, las=1)
158 boxplot(result, medcol="red", col="lightgreen",
159         ylab="Per capita growth, percent")
160 grid(nx=NA, ny=NULL, col="lightgray", lty=1, lwd=1)

162 dev.off()
```

---

**Listing B.2.** Using R with  $\LaTeX$  to create a WDI-based presentation (*extractWDIBeamer.R*).

```

1 # -----
2 #
3 #   Code to extract data series from WDI and create line charts
4 #       for a Latex-based Beamer presentation
5 #
6 #   Original code by: Dominique van der Mensbrugge
7 #
8 # -----
9
10 # >>>> BEGINNING OF USER INPUT
11
12 # Set output directory (use './' for default directory)
13
14 # Windows
15 outwd <- "v:/data/wdi/"
16
17 # MAC
18 # outwd <- "~/Documents/Dominique/data/wdi/"
19
20 # Set file name
21
22 fName <- "pptExample.tex"
23 fName <- paste(outwd, fName, sep="")
24
25 # Set title, author and address (Use \\ to separate lines)
26 pTitle <- "Example of creating slides from an R-based WDI extraction procedure"
27 shortTitle <- "WDI Charts"
28 author <- "Dominique van der Mensbrugge"
29 nickName <- "DvdM"
30 address <- "Center for Global Trade Analysis (GTAP) \\ \\ \\ Purdue University \\ \\ \\
31   West Lafayette, IN"
32 institute <- "GTAP"
33
34 # Set the beginning year
35 begYear <- 1990
36
37 # Set the final year
38
39 endYear <- 2013
40
41 # Select the regions to be extracted
42
43 # The regions are an r x 3 matrix where r is the number of regions
44 # Column 1 is the region ISO-2 code and is used for the extraction
45 # Column 2 is the region ISO-3 code
46 #   (and will be used for the legend, and otherwise can be changed)
47 # Column 3 is the region name (and currently ignored)
48
49 regions <- rbind(cbind("4e", "EAP", "East Asia & Pacific"),
50                 cbind("8S", "SAS", "South Asia"),
51                 cbind("7E", "ECA", "Europe & Central Asia"),
52                 cbind("XQ", "MNA", "Middle East & North Africa"),
53                 cbind("ZF", "SSA", "Sub-Saharan Africa"),
54                 cbind("XJ", "LAC", "Latin America & Caribbean"),
55                 cbind("XD", "HIC", "High-income countries")
56                 )

```

```
58 # Select the indicators

60 # The indicators are an n x 4 matrix were n is the number of indicators
61 # Column 1 is an indicator name (ignored for the moment)
62 # Column 2 is the WB CETS code and is used for the extraction
63 # Column 3 is a scale that will be used to scale the output data
64 # Column 4 describes the indicator -- it will be used as the chart title

66 indTable <- rbind(
67   cbind("pop", "sp.pop.totl", 1e6, "Population (million)"),
68   cbind("gdpkd", "ny.gnp.mktp.kd", 1e9, "GDP (\\$2005 billion)"),
69   cbind("gdpkdpp", "ny.gdp.mktp.pp.kd", 1e12, "GDP (\\$2005PPP trillion)"),
70   cbind("gdppckd", "ny.gnp.pcap.kd", 1e0, "GDP per capita (\\$2005)"),
71   cbind("gdppckdpp", "ny.gdp.pcap.pp.kd", 1e0, "GDP per capita (\\$2005PPP)")
72 )

75 # Set the colors for the lines -- see Latex's 'xcolor' package for choices

77 Colors <- rbind(
78   "blue",
79   "LightSkyBlue",
80   "DarkGrey",
81   "ForestGreen",
82   "red",
83   "black",
84   "Orchid",
85   "Goldenrod",
86   "ForestGreen",
87   "red",
88   "black",
89   "Orchid",
90   "Goldenrod"
91 )

93 # END OF USER INPUT <<<<<

95 # -----
96 #
97 # Main part of WDI extraction
98 #
99 # -----

101 # Load the WDI library

103 library(WDI)

105 # Get the number of indicators and regions

107 nInd <- nrow(indTable)
108 nReg <- nrow(regions)
109 maxColor <- nrow(Colors)

111 # Get the iso2 codes

113 iso2 <- toupper(regions[,1])

115 # Write the Latex file preamble
116 # N.B. It is possible to change 'beamer' themes and theme colors
```

```

118 sink(fName, append=FALSE, split=FALSE)
119 cat("\\documentclass[xcolor={dvipsnames, svgnames}]{beamer}\n")
120 cat("\\usepackage{CambridgeUS}\n")
121 cat("\\usecolortheme{orchid}\n\n")
122 cat("\\usepackage{pgfplots}\n\n")
123 cat("\\begin{document}\n\n")
124 cat("\\pgfplotsset{every axis/.append style={line width=0.75pt},
125 axis x line*=bottom, axis y line*=left}\n")
126 cat("\\pgfsetplotmarksize{0pt}\n\n")
127 cat("\\title[", shortTitle, "]{", pTitle, "}\n", sep="")
128 cat("\\author[", nickName, "]{", author, "}\n", sep="")
129 cat("\\institute[", institute, "]{", address, "}\n", sep="")
130 cat("\\date{\\today}\n")
131 cat("\\frame{\\titlepage}\n\n")
132 sink()

134 # Loop over all indicators

136 for (i in 1:nInd) {

138 # Get the data from the WDI database

140 cube <- WDI(iso2, indicator=indTable[i,2], start=begYear, endYear)

142 # Save the chart pre-amble

144 sink(fName, append=TRUE, split=FALSE)
145 cat("\\frame{\\frametitle{", indTable[i,4], "}}\n", sep="")
146 cat("\\begin{figure}[h]\n")
147 cat("\\centering\n")
148 cat("\\begin{tikzpicture}[scale=0.5]\n")
149 cat("\\begin{axis}[width=19.0cm, height=13.5cm, \n")
150 cat(" xlabel={Year}, \n")
151 cat(" xmin={", begYear-1, "}, \n", sep="")
152 cat(" xmax={", endYear, "}, \n", sep="")
153 cat(" ymajorgrids, \n")
154 cat(" legend style={draw=none}, \n")
155 cat(" x tick label style={/pgf/number format/.cd, scaled x ticks = false, set
thousands separator={}, fixed}, \n")
156 cat(" y tick label style={/pgf/number format/.cd, scaled y ticks = false, set
decimal separator=., fixed}, \n")
157 cat(" legend style={at={(1.03, 0.5)}, anchor=west, font=\\footnotesize}, \n")
158 cat(" legend cell align=left\n")
159 cat(" ]\n")

162 # Set the scale for this indicator and initialize the color table
163 indScale <- as.numeric(indTable[i,3])
164 nColor <- 0

166 # Loop over all regions and write out the data
167 for (r in 1:nReg) {

169 nColor <- nColor + 1
170 if(nColor > maxColor) nColor <- 1

172 # Get the data for the selected region

174 regData <- cube[cube$iso2c==iso2[r], c("iso2c", indTable[i,2], "year")]

```



```
176 # Sort the data by year
177 x <- regData[order(regData$year),]

179 # Get the number of years
180 nYear = nrow(x)

182 # Save the data for this indicator

184 cat("\\addplot+[", Colors[nColor], ", smooth, line width=2.0pt] coordinates\n",
185     sep="")
186 cat("{\n")

187 for(t in 1:nYear) {
188   cat("(" , x[t,3], ", " , x[t,2]/(indScale), ") \n", sep="")
189 }
190 cat("};\n")
191 cat("\\addlegendentry{", regions[r,2], "} \n", sep="")
192 }

194 # End the graph

196 cat("\\end{axis}\n")
197 cat("\\end{tikzpicture}\n")
198 cat("\\end{figure}\n")
199 cat("}\n\n")
200 sink()
201 }

203 # Write the Latex closing

205 sink(fName, append=TRUE, split=FALSE)
206 cat("\\end{document}\n\n")
207 sink()
```

---

**Listing B.3. A GAMS file to retrieve data from WDI (gdpr.gms).**

```
1 *
-----
2 *
3 *   Example of a GAMS file that calls R to extract data from
4 *   the World Bank's WDI database
5 *
6 *   Description:
7 *   This GAMS file provides an example of how to combine GAMS with R to
8 *   extract data from the World Bank's World Development Indicators
9 *   Database (WDI). The user prepares a list of countries and indicators to
10 *   extract from WDI, including the range of years. The GAMS program will
11 *   save all of this information in a text file and then invoke R as a
12 *   sub-process (using R's RScript program). If successful, the R script
13 *   will save the extracted data as a CSV file. The file is converted to
14 *   a GDX file before being read back in to the GAMS file. This is
15 *   necessary because CSV files are read-in at compile time, but GDX files
16 *   are read in at execution time. A utility called CSV_GDX_Tools is used
17 *   to convert the R-outputted CSV file, to the GDX format.
18 *   N.B. When reading CSV files GAMS does not want a header line in the
19 *   data file that describes the different fields in the CSV file.
20 *   However the CSV to GDX converter, and many database-type translators
21 *   of CSV files, such as Excel's pivot table feature require
22 *   the header line. This is one of the user options provided below.
23 *
24 *   Requirements:
25 *
26 *   isocodes.gms: A file containing a current set of ISO-2 and ISO-3 codes
27 *   with their mapping. In addition, the file contains a
28 *   mapping from ISO-3 to the GTAP regions.
29 *
30 *   R:           The RScript function is used to run R in background.
31 *               When invoked, RScript should be in the system's path.
32 *               If not, invoke R with full path name.
33 *
34 *   CSV_GDX_Tools The CSV to GDX conversion program. It needs to be
35 *               in the system's path, or else, the full path name
36 *               needs to be provided. (N.B. The tool also requires
37 *               a file called 'forbiddenwords.txt' that should be
38 *               in the same folder as the conversion tool.)
39 *
40 *   extractWDI.R The name of the R script file that performs the
41 *               extraction from WDI and saves the data as a CSV file.
42 *               Users are welcome to modify and/or rename the file.
43 *
44 *   User inputs:
45 *
46 *   BASENAME     This will be the base file name for all of the
47 *               files created by the program.
48 *
49 *   FIRSTYR      First year to extract. The extraction routine requires a
50 *               continuous range.
51 *
52 *   LASTYR       The last year to extract.
53 *
54 *   EXTALL       A true or false flag. If 'true' the extraction routine
55 *               will download the data for all countries in the WDI
56 *               database for the selected indicators. Any country
57 *               information (for example as defined in the subset 'c2'
```

```

58 *           will be ignored. If 'false', only data for countries
59 *           identified in the subset 'c2' will be extracted.
60 *
61 *   IFHEADER   A true or false flag. If 'true' a header line will be
62 *             included at the top of the CSV file with labels for the
63 *             CSV file fields. If 'false', no header line will be
64 *             output. For reading the CSV data into GAMS, the flag
65 *             should be set to 'false'.
66 *
67 *   IFNAME     A true or false flag. The WDI extraction also retrieves
68 *             the country names. If this flag is set to 'false', the
69 *             country names will not be saved with the data. These can
70 *             nonetheless be retrieved by the GAMS code since the
71 *             country names are provided with the set definitions. The
72 *             information is therefore redundant and takes up storage
73 *             space (and can also not be readily read into a GAMS
74 *             matrix). Set to 'true' to save the country names with
75 *             the CSV data.
76 *
77 *   c2         C2 is a subset of the ISO-2 codes (that are used by the
78 *             WDI extraction routine). The subset is user-defined and
79 *             contains the ISO-2 codes for the countries which are to
80 *             be extracted. The subset is ignored if 'EXTALL' is set
81 *             to 'true'.
82 *
83 *   var        Var is a set of indicators to extract. The user provides
84 *             short names for the extracted indicators--though is free
85 *             to use the same mnemonic as the World Bank database. The
86 *             set Var must be paired with the set wbind and the mapv
87 *             mapping (see below).
88 *
89 *   wbind      WBInd is a set of indicators to extract using the WB's
90 *             mnemonics. It is typically paired with the set Var and
91 *             the mapping mapv.
92 *
93 *   mapv       The set mapv pairs the user-named indicators (Var) with
94 *             the WB-named indicators (WBInd). It should be a
95 *             one-to-one mapping.
96 *
97 *
-----
99 acronym true, false ;

101 * Read the ISO-2, ISO-3 and GTAP codes and mappings

103 $offlisting
104 $include "isocodes.gms"
105 $onlisting

107 *
-----
108 *
109 * BEGINNNING OF USER INPUT
110 *
111 *
-----

113 * System options
114 $setglobal R_EXE C:\Program Files\R\R-3.2.2\bin\RScript

```

```

115 $setglobal CSVGDX "V:\bin\CSV_GDX_Tools"
116 $setglobal RSCRIPT "extractWDI.R"

118 * User options
119 $setglobal BASENAME "GDPR"
120 $setglobal FIRSTYR 2000
121 $setglobal LASTYR 2014
122 $setglobal EXTALL false
123 $setglobal IFHEADER true
124 $setglobal IFNAME false

126 * NOTE: If EXTALL is set to true, the following subset will be ignored
127 *       no matter how it is defined.

129 set c2(iso2) "Countries to extract" / fr, us, gb, be, bz / ;

131 set var "Indicators to extract" /
132     "pop"
133     "gdpcd"
134     "gdpkd"
135     "gnpppcd"
136     "gnpppkd"
137 / ;

139 set wbind "WB indicators to extract" /
140     "SP.POP.TOTL"
141     "NY.GDP.MKTP.CD"
142     "NY.GDP.MKTP.KD"
143     "NY.GDP.PCAP.PP.CD"
144     "NY.GDP.PCAP.PP.KD"
145 / ;

147 set mapv(var, wbind) "Mapping of variable names to WB indicators" /
148     "pop" . "SP.POP.TOTL"
149     "gdpcd" . "NY.GDP.MKTP.CD"
150     "gdpkd" . "NY.GDP.MKTP.KD"
151     "gnpppcd" . "NY.GDP.PCAP.PP.CD"
152     "gnpppkd" . "NY.GDP.PCAP.PP.KD"
153 / ;

155 parameter scale(var) "Scaling factor for indicators" /
156     "pop" 1e-6
157     "gdpcd" 1e-6
158     "gdpkd" 1e-6
159     "gnpppcd" 1e-0
160     "gnpppkd" 1e-0
161 / ;

163 *
-----
164 *
165 * END OF USER INPUT
166 *
167 *
-----

169 sets
170 year "Years to extract" / %FIRSTYR%*%LASTYR% /
171 y0(year) "First extraction year" / %FIRSTYR% /
172 yf(year) "Last extraction year" / %LASTYR% /

```

```

173 ;

175 set c(iso3) "Countries being extracted with ISO-3 codes" ;

177 if(%EXTALL% eq false,
178     loop(mapISO(iso2,iso3)$c2(iso2),
179         c(iso3) = yes ;
180     ) ;
181 else
182     c(iso3) = yes ;
183 ) ;

185 * Delete intermediate files if they exist

187 $if exist "%BASENAME%.opt" $call 'del "%BASENAME%.opt"'
188 $if exist "%BASENAME%.csv" $call 'del "%BASENAME%.csv"'

190 file fopt / "%BASENAME%.opt" / ;
191 put fopt ;
192 put '"Name","Value"' / ;

194 put '"Output file", "%BASENAME%.csv"' / ;

196 loop(y0, put '"Start year", ', y0.tl / ; ) ;
197 loop(yf, put '"End year", ', yf.tl / ; ) ;

199 if(%IFHEADER% eq false,
200     put '"Header", "FALSE"' / ;
201 else
202     put '"Header", "TRUE"' / ;
203 ) ;

205 if(%IFNAME% eq false,
206     put '"Country name", "FALSE"' / ;
207 else
208     put '"Country name", "TRUE"' / ;
209 ) ;

211 scalar ifFirst / 1 / ;

213 put '"Regions", "' ;
214 if(%EXTALL% eq false,
215     loop(c2,
216         if(ifFirst,
217             put c2.tl:2 ;
218             ifFirst = 0 ;
219         else
220             put ', ', c2.tl:2 ;
221         ) ;
222     ) ;
223 else
224     put "all" ;
225 ) ;

227 put '"'" / ;
228 put '"Indicators", "wbName"' / ;
229 fopt.pc=5 ;
230 loop(mapv(var,wbind),
231     put var.tl, wbind.tl / ;
232 ) ;

```

```
233 putclose fopt ;

235 * Extract the data using R
236 execute '%R_EXE%' %RSCRIPT% %BASENAME%.opt'

238 * Convert the data to GDX format
239 execute '%CSVGDX%' "%BASENAME%.csv" comma "All" "Var,iso3,Year" /method=csvgdx /
  PARNAME=wdiData /GDX=%BASENAME%.gdx'

241 * Load the GDX data

243 parameter
244   wdiData(var,iso3,year)
245 ;
246 execute_load "%BASENAME%.gdx", wdiData ;

248 * Process the read-in data
249 wdiData(var,c,year) = scale(var)*wdiData(var,c,year) ;
250 option decimals=0 ;
251 display wdiData ;

253 * Cleanup intermediate files
254 if(1,
255 $if exist "%BASENAME%.opt" $call 'del "%BASENAME%.opt"'
256 $if exist "%BASENAME%.csv" $call 'del "%BASENAME%.csv"'
257 $if exist "%BASENAME%.gref" $call 'del "%BASENAME%.gref"'
258 ) ;
```

---

**Listing B.4.** WDI.Retrieve function (*extractWDI.R*).

---

```
1 wdi.Retrieve <- function(cName, sName, indName, startYear, endYear,
2   fName, ifAppend, ifName) {
3   #
4   # Arguments:
5   #   cName: List of countries/regions being extracted (could be "all")
6   #   sName: Short series name (will be output to CSV file)
7   #   indName: Mnemonic of WDI indicator
8   #   startYear: Beginning year
9   #   endYear: End year
10  #   fName: Name of output file
11  #   ifAppend: FALSE for first record (normally header), otherwise TRUE
12  #   ifName: FALSE to delete country name from output, otherwise TRUE
13  #
14
15  # Retrieve the indicator
16
17  cube <- WDI(country=cName, indicator=indName, start=startYear, end=endYear)
18
19  # Convert the ISO-2 labels to ISO-3 labels -- requires the countrycode package
20
21  cube$iso3c <- countrycode(cube$iso2c, "iso2c", "iso3c")
22
23  # Fix ISO-3 codes for missing countries
24
25  iso2Ndx <- grep('iso2c', colnames(cube))
26  iso3Ndx <- grep('iso3c', colnames(cube))
27
28  # !!!! May need to add other countries and/or change ISO-3 codes
29
30  cube[iso3Ndx][cube[iso2Ndx]=='JG'] <- "JGY"
31  cube[iso3Ndx][cube[iso2Ndx]=='XK'] <- "KOS"
32
33  # Delete regions with no iso code
34
35  cube <- cube[!(is.na(cube$iso3c)),]
36
37  # Append the variable name to the dataframe
38
39  cube[, "Var"] <- sName
40
41  # Create the final data frame by re-ordering columns
42
43  if(ifName) {
44
45    cube <- cube[,c("Var", "iso3c", "country", "year", indName)]
46
47  }
48  else {
49
50    cube <- cube[,c("Var", "iso3c", "year", indName)]
51
52  }
53
54  # Save the data in CSV format
55
56  write.table(cube, fName, append=ifAppend, row.names=FALSE,
57    col.names=FALSE, na="", sep=",")
58 }
```

---

**Listing B.5.** Main part of R WDI extraction script used with GAMS (*extractWDI.R*).

---

```
2 #
-----
3 #
4 # Main part of WDI extraction
5 #
6 #
-----

8 # Load the needed libraries (after installation on the local machine)

10 library(WDI)
11 library(countrycode)

13 # Read the command line arguments

15 args <- commandArgs(TRUE)

17 if(length(args) < 1) {
18   srcFile <- "extractWDI.opt"
19 } else {
20   srcFile <- args[1]
21 }

23 # Get the options from the options file

25 options <- read.csv(srcFile, stringsAsFactors = FALSE,
26   strip.white=TRUE, quote = "\"\"")

28 # Get the name of the output file

30 outFile <- options[grep("Output file", options$Name),2]

32 # Get the beginning year

34 begYear <- as.numeric(options[grep("Start year", options$Name),2])

36 # Get the final year

38 endYear <- as.numeric(options[grep("End year", options$Name),2])

40 # Flag for header

42 ifHeader <- as.logical(options[grep("Header", options$Name),2])

44 # Flag for outputting country name

46 ifName <- as.logical(options[grep("Country name", options$Name),2])

48 # Get the regions -- 'all' is to get all regions, must use iso2c codes

50 regions <- options[grep("Regions", options$Name),2]

52 if(regions != "all") {
53   regions <- toupper(regions)
54 }

56 # Parse the regions
```



```
58 regions <- strsplit(regions, ",")
59 regions <- c(do.call("cbind",regions))

61 # Get the indicators -- assumes input file has no blank lines

63 row1 <- 8
64 row2 <- nrow(options)
65 indTable <- options[row1:row2,]

67 # Get the number of indicators

69 nInd <- nrow(indTable)

71 if(ifHeader) {

73 # Use header for CSV file
74   if(ifName) {
75     header <- cbind("Var", "iso3", "regName","Year", "Val")
76   }
77   else {
78     header <- cbind("Var", "iso3", "Year", "Val")
79   }

81   write.table(header, outFileNames, append=FALSE, row.names=FALSE,
82             col.names=FALSE, sep=",")
83 }
84 ifAppend <- ifHeader

86 # Loop over all indicators

88 for (i in 1:nInd) {

90 # Get the short and World Bank name for the indicators

92   var <- indTable[i,1]
93   wbName <- indTable[i,2]

95 # Use the function above to retrieve the data, convert iso2 to iso3, drop the
96 #   regions, and append to the CSV file

98   wdi.Retrieve(regions, var, wbName, begYear, endYear,
99               outFileNames, ifAppend, ifName)

101 # After the first indicator, ifAppend must be TRUE

103   ifAppend <- TRUE
104 }
```

---