

UN Sustainable Development Goals Data Analysis

Fredrik Randsborg Bølstad

June 16, 2019

Contents

1	Introduction	2
2	Data analysis	2
2.1	Downloading data	2
2.1.1	UN SDG Indicators data from Google BigQuery	2
2.1.2	Gapminder country list	3
2.1.3	World Bank population data	4
2.2	Cleaning data	4
2.2.1	Merging datasets	4
2.2.2	Selecting relevant data	5
2.2.3	Dropping, renaming and coercing data	6
2.2.4	Calculate absolute population numbers	6
2.2.5	Append target descriptions for SDG indicators	7
2.2.6	Harmonizing map data with the main dataset	7
2.2.7	Save tidy data to disk	7
2.3	Exploring and visualizing data	8
2.3.1	World heatmaps	8
2.3.2	World heatmap - rural vs urban grid	14
2.3.3	Scatterplot - rural vs urban grid	15
2.3.4	Combined heatmap and scatterplot - rural vs urban grid	17
2.3.5	Animated heatmaps and scatterplots	18
2.3.6	Country bar chart comparison	18
2.3.7	Wide form time-series table	21
2.4	Forecasting	24
2.4.1	Fitting a linear model using <code>lm</code>	25
3	Results	27
4	Conclusion	27

1 Introduction

The purpose of this project is to explore the global progress towards the United Nations (UN) Sustainable Development Goals (SDG). The main dataset is the UN SDG Indicators dataset provided by the UN Statistics Division. It is complimented by data from the Gapminder Foundation and the World Bank. Finding powerful ways of visualizing the dataset will be of particular importance to the project. With the help of the maps and mapdata R packages, a key ambition is to visualize reported progress for each country on a world map. Finally, the report will touch on what the current data may say about future results. Key goals:

- Programmatically download data from the relevant sources
- Clean and merge data into unified datasets
- Communicate reported results graphically
- Use trend analysis to predict future results

2 Data analysis

2.1 Downloading data

The main data for this project is provided by the UN Statistics Division via Google BigQuery. The geographic areas in the UN dataset includes countries, regions and continents. The geographical scope of this analysis is limited to the country list used by the Gapminder Foundation. The UN dataset will thus have to be filtered to only keep the countries that are in the Gapminder's data. Population data is separately downloaded from The World Bank's repository. Country coordinate data is extracted from the maps and mapdata R packages.

Table 1: Data sources

Data source	Fetch method
UN Sustainable Development Goals Indicators	bigrquery - R interface to BigQuery
Gapminder geography data	googlesheets
World Bank population data	wbstats - R interface to World Bank API
Spatial map data	maps and mapdata

The raw data is downloaded and saved for processing. The raw and clean data are thus stored separately to avoid having to download the raw data over again, in case an unexpected or undesired result flow from modification of the dataset during data cleaning.

The code used to download and save the above data from each source can be found below. The code is commented to explain exactly what it does.

2.1.1 UN SDG Indicators data from Google BigQuery

```
library(bigrquery)      # Query Google's BigQuery database
library(DBI)            # Interface for communication with relational database systems

# Local destination for downloaded raw data
un_sdg_file = "./raw_data/bigquery_un_sdg.rda"

# Google Cloud project ID required for BigQuery interactions
```

```

# Create your own project ID at https://console.cloud.google.com/
google_cloud_project_id <- "edx-ds-capstone-2019"

# BigQuery connection interface
dbi_connection <- dbConnect(
  bigrquery::bigrquery(),
  project = "publicdata",
  dataset = "un_sdg",
  billing = google_cloud_project_id
)

# SQL query for entire dataset (~300 MB)
sql_query <- "SELECT * FROM `bigquery-public-data.un_sdg.indicators`"

# Download dataset
un_sdg_data <- dbGetQuery(dbi_connection, sql_query)

# Save dataframe to file
save(un_sdg_data, file = un_sdg_file)

```

2.1.2 Gapminder country list

```

library(googlesheets)    # Google spreadsheet API

# Local destination for downloaded raw data
gapminder_file = "./raw_data/gapminder_country_list.rda"

# Connect to spreadsheet (URL from https://www.gapminder.org/data/geo/)
gapminder_spreadsheet <- gs_url(paste("https://docs.google.com/spreadsheets/",
                                       "d/1qHalit8sXC0R8oVXibc2wa2gY7bkWGzOybEMTWp-08o/",
                                       sep = ""))
# formatted to fit on page

# Download list-of-countries worksheet
gapminder_country_list <- gapminder_spreadsheet %>%
  gs_read(ws = "list-of-countries-etc")

# Keep only country and region.
#> Mutate colname to 'geoareaname' to match UN SDG data colname
#> Make geo uppercase to match World Bank data
gapminder_country_list <- gapminder_country_list %>%
  mutate(geoareaname = name, region = four_regions) %>%
  select(geo, geoareaname, region) %>%
  mutate(geo = toupper(geo))

# Save dataframe to file
save(gapminder_country_list, file = gapminder_file)

```

2.1.3 World Bank population data

```
library(wbstats) # World Bank API interface

# Local destination for downloaded raw data
wb_pop_file = "./raw_data/world_bank_population_data.rda"

# Create list of ISO 3 country codes from Gapminder
# which will be used to query World Bank's data
iso3c_codes <- gapminder_country_list$geo

# Remove HOS (Holy See) which isn't in World Bank data
# https://datahelpdesk.worldbank.org/knowledgebase/articles/898590-country-api-queries
iso3c_codes <- iso3c_codes[-which(iso3c_codes == 'HOS', )]

# Get year interval for which there is UN SDG data
if(!exists("un_sdg_data")) load(un_sdg_file);
first_year <- min(un_sdg_data$timeperiod)
last_year <- max(un_sdg_data$timeperiod)

# Download population size from the World Bank
wb_pop <- wb(country = iso3c_codes,
  indicator = 'SP.POP.TOTL',
  startdate = first_year, enddate = last_year)

# Save dataframe to file
save(wb_pop, file = wb_pop_file)
```

2.2 Cleaning data

2.2.1 Merging datasets

The data from the UN, Gapminder Foundation and the World Bank are merged into one unified dataset. There are two main goals for this process:

1. Join UN SDG data with Gapminder data, such that all rows relating to geographical areas not in the Gapminder country list are dropped.
2. Append the World Bank population data for each country and the respective years in the dataset.

```
### Join UN SDG and Gapminder dataset such that:
# > Drop from UN SDG all countries not in Gapminder
# > Append to UN SDG the region column from Gapminder
dat <- inner_join(un_sdg_data, gapminder_country_list)

### Append World Bank population data, by country and year
# Rename and select World Bank columns to match UN SDG dataset
wb_pop <- wb_pop %>%
  rename(geo = iso3c, timeperiod = date, population = value) %>%
  select(geo, timeperiod, population)

# Append population column to UN SDG dataset
dat <- left_join(dat, wb_pop)
```

2.2.2 Selecting relevant data

After preliminary exploration of the dataset, including manual inspection, it became clear that the reported UN data were severely lacking. Many indicators had no, little or inconsistent data reporting. Only a minority of indicators were consistently or semi-consistently reported on by a broad community of countries. A selection of indicators with sufficient data were therefore handpicked to be used for the continued analysis. The target description for each indicator were manually added to the dataset for convenience (so that the goal of the indicator can quickly be queried). The complete description of all indicators in the dataset can be found in the annex of the UN General Assembly Resolution [A/RES/71/313](#). This is also the source of the text descriptions copied into the code below.

Table 2: Selected UN SDG Indicators

Indicator	Description
6.1.1	Population proportion using safe drinking water services
7.1.1	Population proportion with access to electricity
7.2.1	Renewable energy (as share of total consumption)
8.1.1	Annual growth rate, real GDP per capita
9.4.1	CO ₂ emission per unit of value added
12.2.1	Material footprint
12.2.2	Domestic material consumption
17.6.2	Population proportion with broadband subscriptions
17.8.1	Population proportion using the internet

```
### Indicators to keep
indicators <- c("6.1.1", "7.1.1", "7.2.1", "8.1.1", "9.4.1",
             "12.2.1", "12.2.2", "17.6.2", "17.8.1")

### Text description of targets for indicators (https://undocs.org/A/RES/71/313)
target_description <-
  c('By 2030, achieve universal and equitable access to
    safe and affordable drinking water for all',
    'By 2030, ensure universal access to affordable, reliable and modern energy services',
    'By 2030, increase substantially the share of
      renewable energy in the global energy mix',
    'Sustain per capita economic growth in accordance with
      national circumstances and, in particular, at least 7 per cent
      gross domestic product growth per annum in the least developed countries',
    'By 2030, upgrade infrastructure and retrofit industries to make them sustainable,
      with increased resource-use efficiency and greater adoption of clean and
      environmentally sound technologies and industrial processes, with all countries
      taking action in accordance with their respective capabilities',
    'By 2030, achieve the sustainable management and efficient use of natural resources',
    'By 2030, achieve the sustainable management and efficient use of natural resources',
    'Enhance North-South, South-South and triangular regional and international
      cooperation on and access to science, technology and innovation and
      enhance knowledge-sharing on mutually agreed terms, including through
      improved coordination among existing mechanisms, in particular at the
      United Nations level, and through a global technology facilitation mechanism',
    'Fully operationalize the technology bank and science, technology and innovation
      capacity-building mechanism for least developed countries by 2017 and enhance the
      use of enabling technology, in particular information and communications technology')
```

```
# Data frame of indicators and text descriptions (later appended to the main dataset)
descriptions_dat <- data.frame(indicators, target_description) %>%
  rename(indicator = indicators)
```

2.2.3 Dropping, renaming and coercing data

The entire UN SDG Indicators dataset were downloaded to begin with. After the preliminary exploration of the data and the selection of indicators, many variables are irrelevant to this project. Those are dropped from the data frame that will later be used for visualization and analysis. A few columns are also renamed to something more descriptive (e.g. *geoareaname* has been filtered to only include countries and is thus renamed to *country*). Finally, each column is coerced into either a numeric (e.g. reported values), integer (e.g. years) or factor (e.g. region) to optimize performance.

```
### Manipulate columns with dplyr
dat <- dat %>%

# Filter for chosen indicators
filter(indicator %in% indicators) %>%

# Remove seriescode, geoareacode, time_detail, freq
select(-c('goal', 'seriescode', 'geoareacode', 'time_detail', 'freq')) %>%

# Rename geoareaname to country, timeperiod to year
rename(country = geoareaname, year = timeperiod) %>%

# Capitalize region column
mutate(region = sub("(.)", "\U001", region, perl=TRUE)) %>%

# Remove empty columns
select_if(~sum(!is.na(.)) > 0) %>%

# Coerce the type of each column
mutate(target = as.factor(target),
       indicator = as.factor(indicator),
       year = as.integer(year),
       value = as.numeric(value),
       nature = as.factor(nature),
       location = as.factor(location),
       type_of_product = as.factor(type_of_product),
       type_of_speed = as.factor(type_of_speed),
       units = as.factor(units),
       region = as.factor(region))
```

2.2.4 Calculate absolute population numbers

The UN SDG data only lists a value for each indicator, e.g. the population proportion that has access to safe drinking water. With the population data from the World Bank, it is also possible to estimate the absolute number of people with access to safe drinking water. This is interesting in itself and could be relevant for future analysis. Future analysis could e.g. weigh global progress according to population numbers instead of treating each country as equally important to aggregate global progress.

```
# Use value and population to create column containing the absolute number of people
dat <- dat %>%
  mutate(nr_people = as.integer(floor(population*value/100)))
```

2.2.5 Append target descriptions for SDG indicators

The target descriptions for each indicator, which were previously saved in a separate data frame, is appended to the main data frame.

```
# Append target descriptions to dat
dat <- left_join(dat, descriptions_dat)
```

2.2.6 Harmonizing map data with the main dataset

Map data from the *map_data* R package is saved to a variable and the country column is renamed to match with the main dataset.

```
### Save map coordinates in dataframe
map_world <- map_data('world')

# Rename region as country to harmonize with 'dat'
map_world <- map_world %>%
  rename(country = region)
```

Find and list all instances of spelling discrepancies for country names.

```
# List all mismatches in country names between dat and map_world
dat %>%
  anti_join(map_world, by = 'country') %>%
  select(country) %>%
  unique()
```

Rename countries for which there are spelling discrepancies in order to cooperate fully with the main dataset.

```
# Rename countries in dat to match those of map_world
dat <- dat %>%
  mutate(country = recode(country,
    'Antigua and Barbuda' = 'Antigua',
    'Holy See' = 'Vatican',
    'Trinidad and Tobago' = 'Trinidad')) %>%

# Remove rows with country Tuvalu
filter(country != 'Tuvalu')
```

2.2.7 Save tidy data to disk

```
### Save data frame to file
save(dat, file = tidy_file)
save(map_world, file = world_map_file)
```

2.3 Exploring and visualizing data

Most of the work on this data analysis project ended up being dedicated to visualizing the data. A significant amount of thought and research went into different options. A key achievement of this project was the successful creation of a function that generates a progress heatmap of the world, given any UN SDG indicator and year. Moreover, data for population proportion indicators were separated into the rural and urban population; a side-by-side comparison of the world heatmaps for each population category, and with a scatterplot showing progress per year, was the culmination of this visualization effort. Several functions were programmed to achieve this result in an intuitive and user friendly manner.

Functions for generating GIFs of these maps were also created. This enables a year-by-year animation of progress. Because the GIFs can't be displayed in this PDF report, the code is not included here. Examples of GIF animations are however uploaded to this project's github repository and linked to in the appropriate section below. The full code for the GIF creation is in the main code file (also available on github).

Titles, subtitles, units, year and regions are programmatically added to the plots.

2.3.1 World heatmaps

The functions needed to generate a world heatmap for any indicator:

```
### Function to filter for given indicator and year; delete empty columns
select_ind_year <- function(data, ind, yr) {

  data %>%
    filter(indicator == ind, year %in% yr) %>%
    select_if(~sum(!is.na(.)) > 0)

}

### Function to create world heatmap
plot_world_heatmap <- function(data, ind, year) {

  # Create filtered dataset
  data_filtered <- select_ind_year(data, ind, year)

  # Join filtered dataset with world map
  map_data <- left_join(map_world, data_filtered, by = 'country')

  # Text for title and legend
  title <- map_data$seriesdescription %>% na.omit() %>% unique()
  unit <- map_data$units %>% na.omit() %>% unique() %>% as.character()

  # Plot graph
  ggplot(map_data, aes(x = long, y = lat, group = group)) +
    geom_polygon(aes(fill = value)) +
    # Titles and axis
    labs(title = paste(title),
        subtitle = paste('Year', year),
        fill = paste(unit),
        x = NULL, y = NULL) +
    # Legend
```

```

scale_fill_gradient(low = "#56B1F7", high = "#132B43", na.value = 'lightgrey') +
# Theme adjustments
theme(plot.title = element_text(size = 9, face = 'bold'),
plot.subtitle = element_text(size = 8, hjust = 0),
axis.ticks = element_blank(),
axis.text = element_blank(),
legend.key.size = unit(14, 'pt'),
legend.title = element_text(size = 7),
legend.text = element_text(size = 7),
panel.background = element_rect(fill = 'white'),
plot.background = element_rect(fill = 'white'),
panel.border = element_rect(size = 1, colour = 'gray50', fill = NA))

}

```

Looping through the above function for all selected indicators:

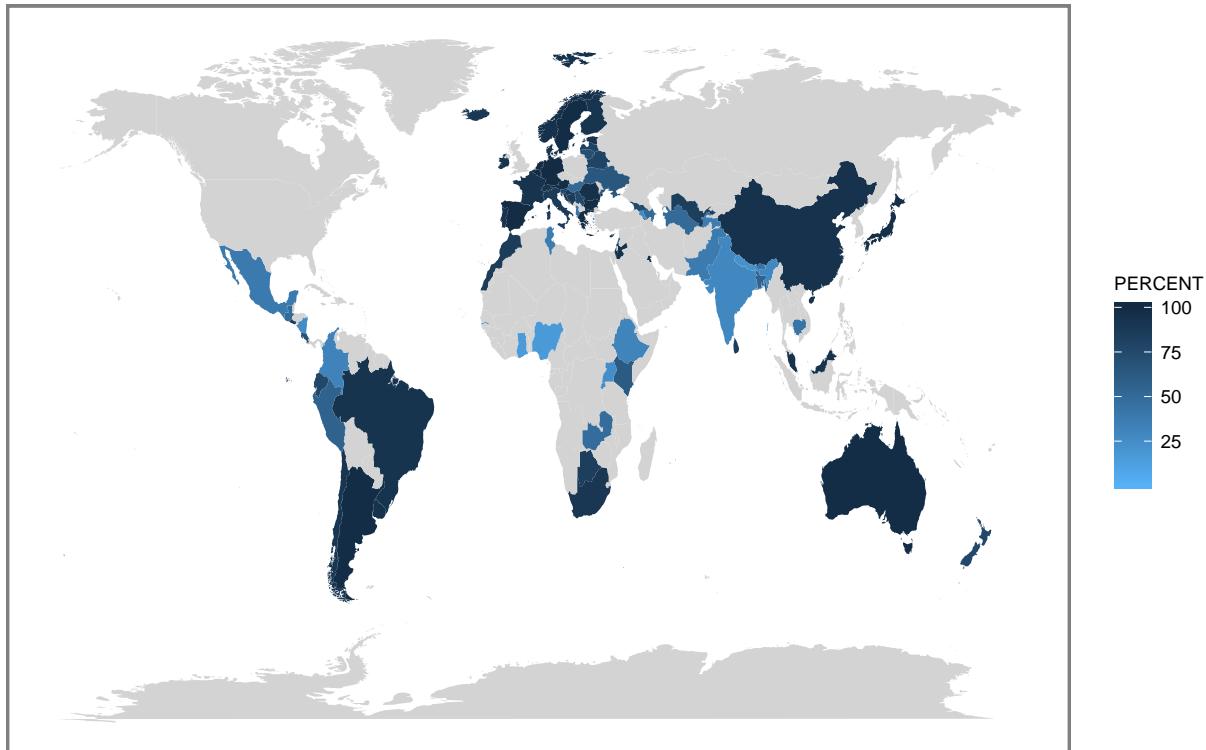
```

for (ind in indicators) {
  heatmap <- plot_world_heatmap(dat, ind, 2000)
  print(heatmap)
}

```

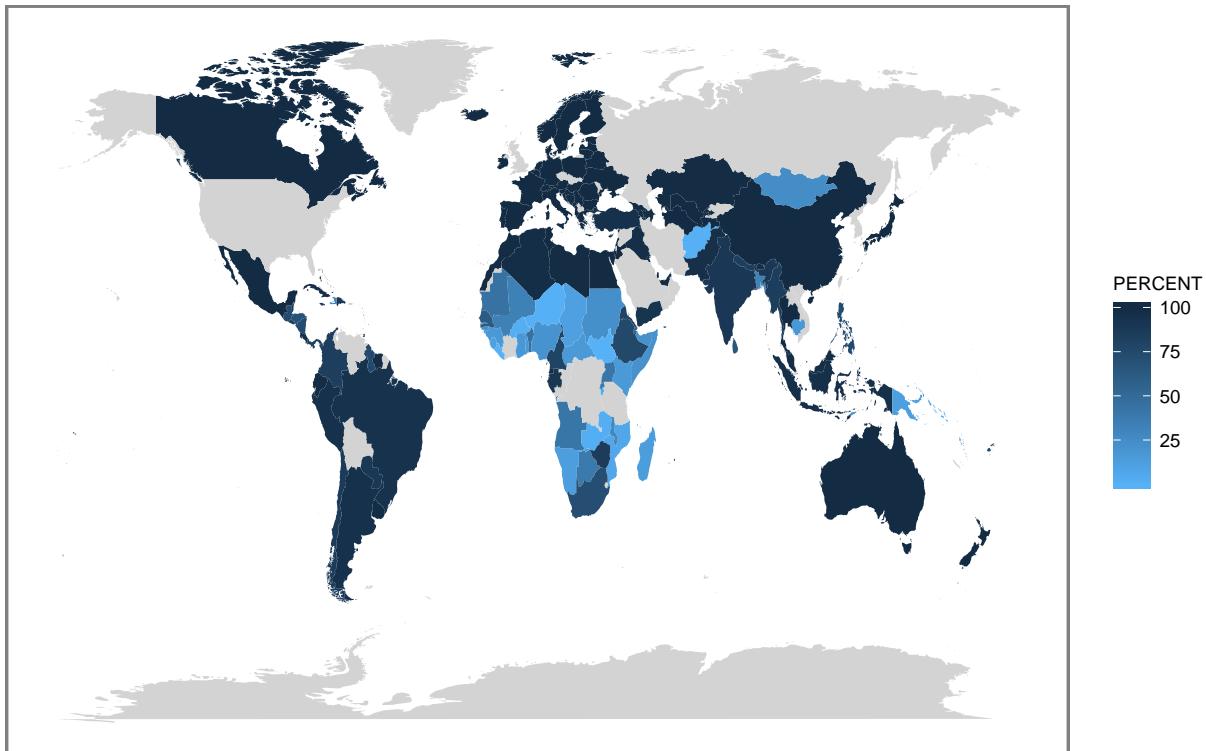
Proportion of population using safely managed drinking water services, by urban/rural (%)

Year 2000



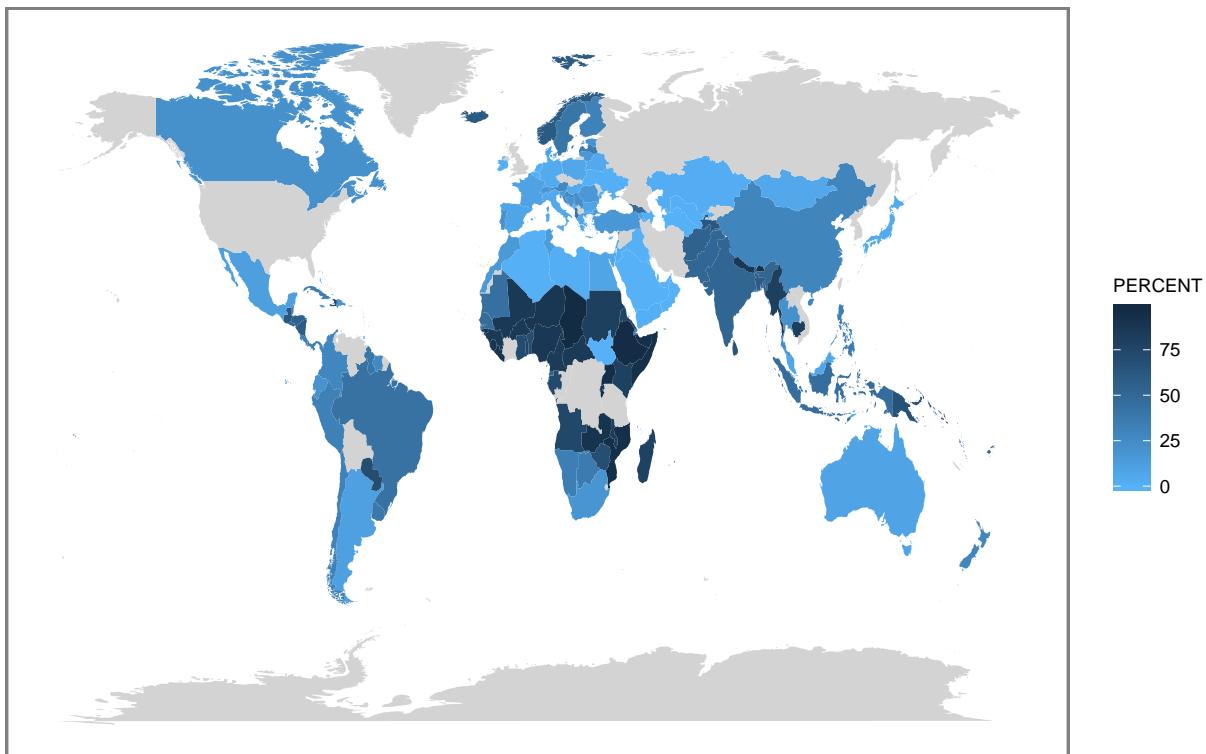
Proportion of population with access to electricity, by urban/rural (%)

Year 2000



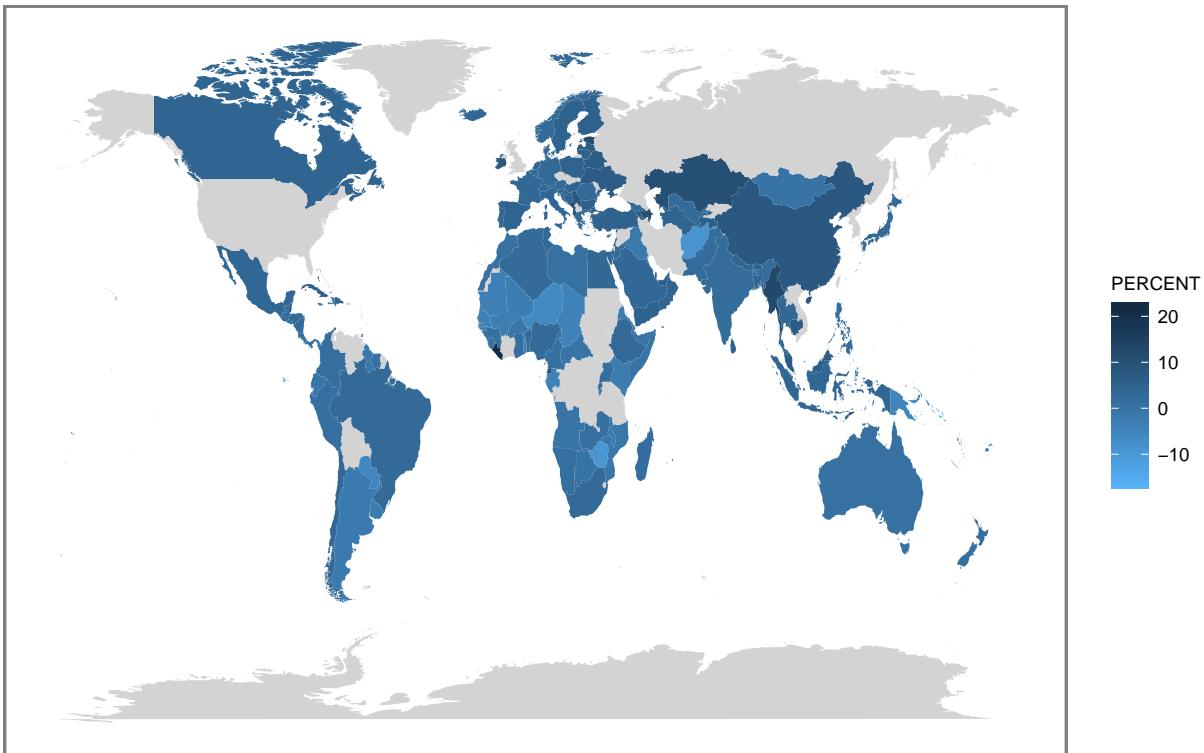
Renewable energy share in the total final energy consumption (%)

Year 2000



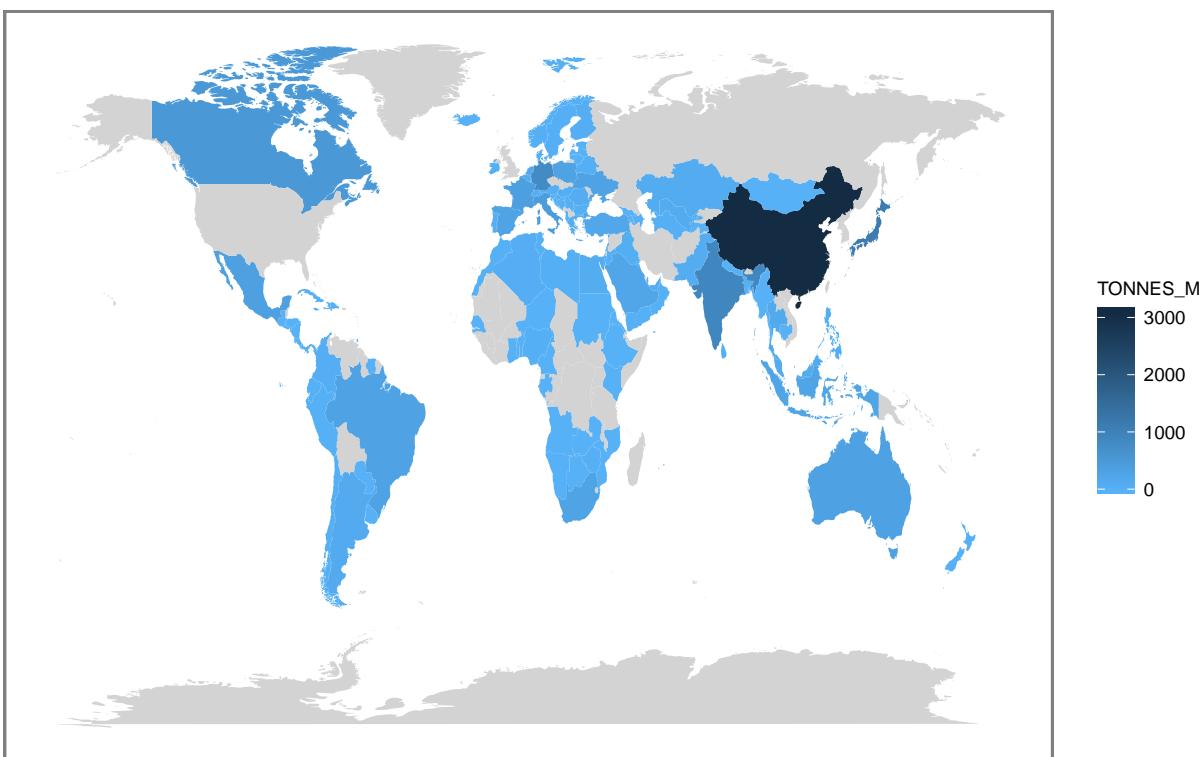
Annual growth rate of real GDP per capita (%)

Year 2000



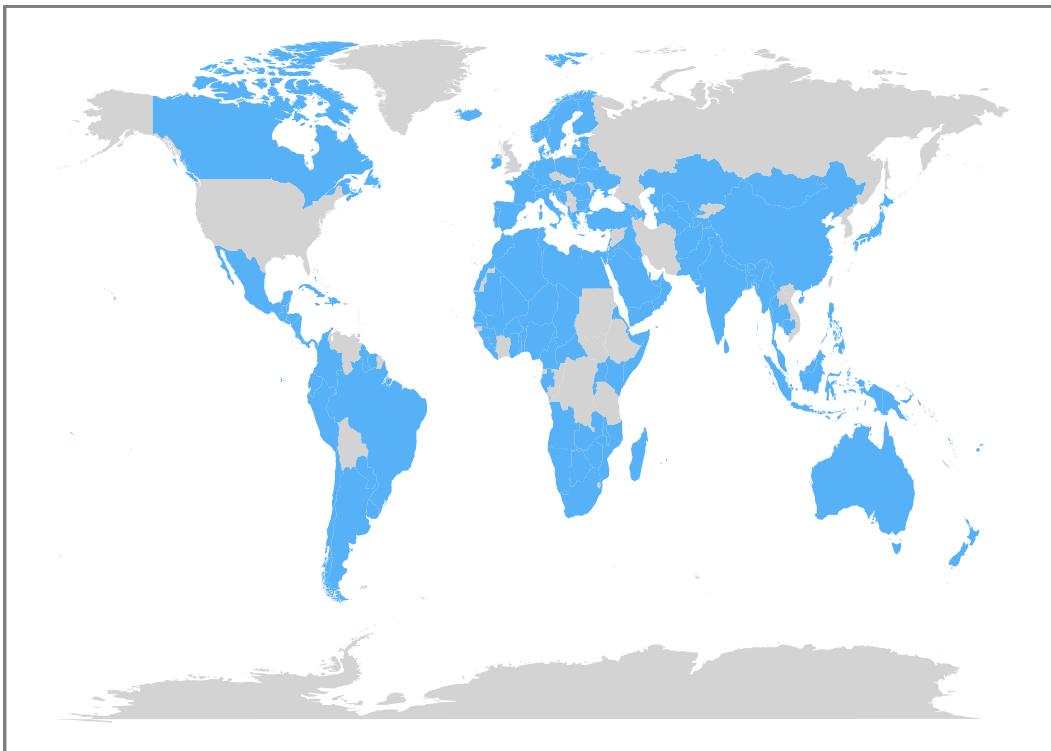
Carbon dioxide emissions from fuel combustion (millions of tonnes)

Year 2000



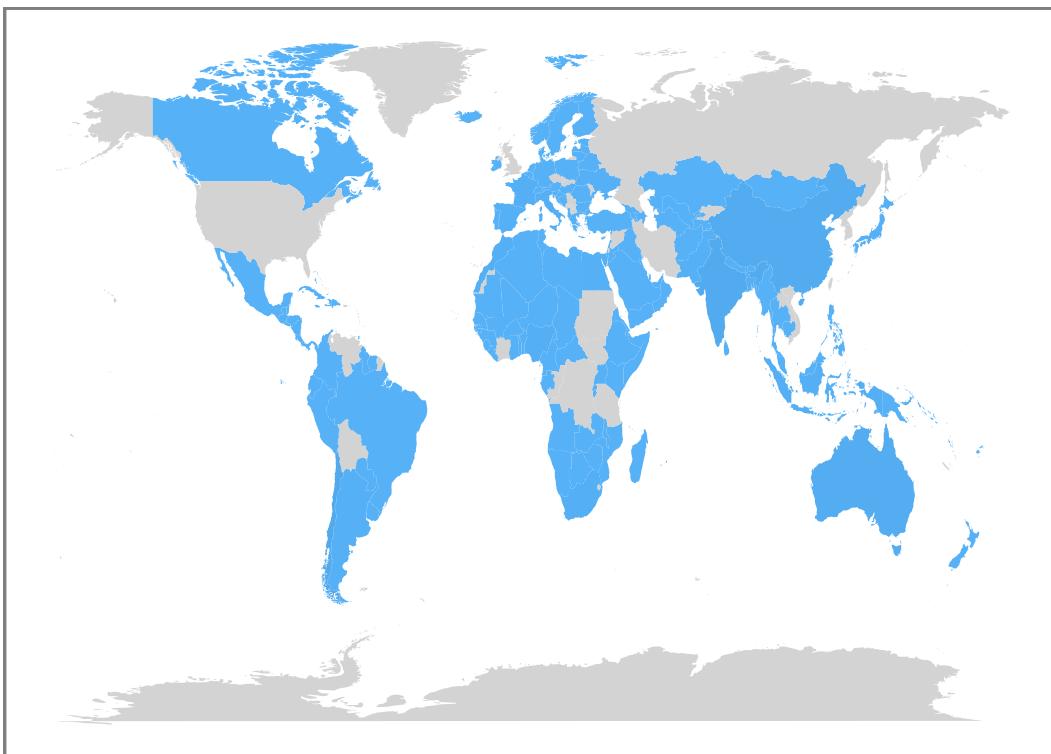
Material footprint per capita, by type of raw material (tonnes)

Year 2000



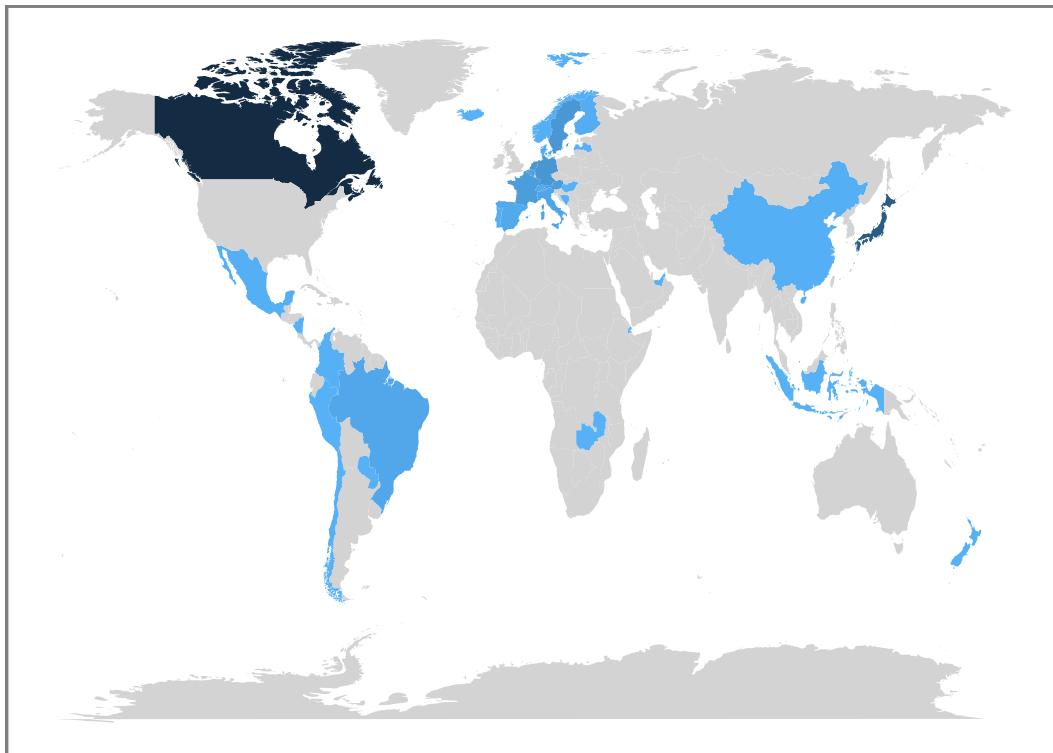
Domestic material consumption per capita, by type of raw material (tonnes)

Year 2000



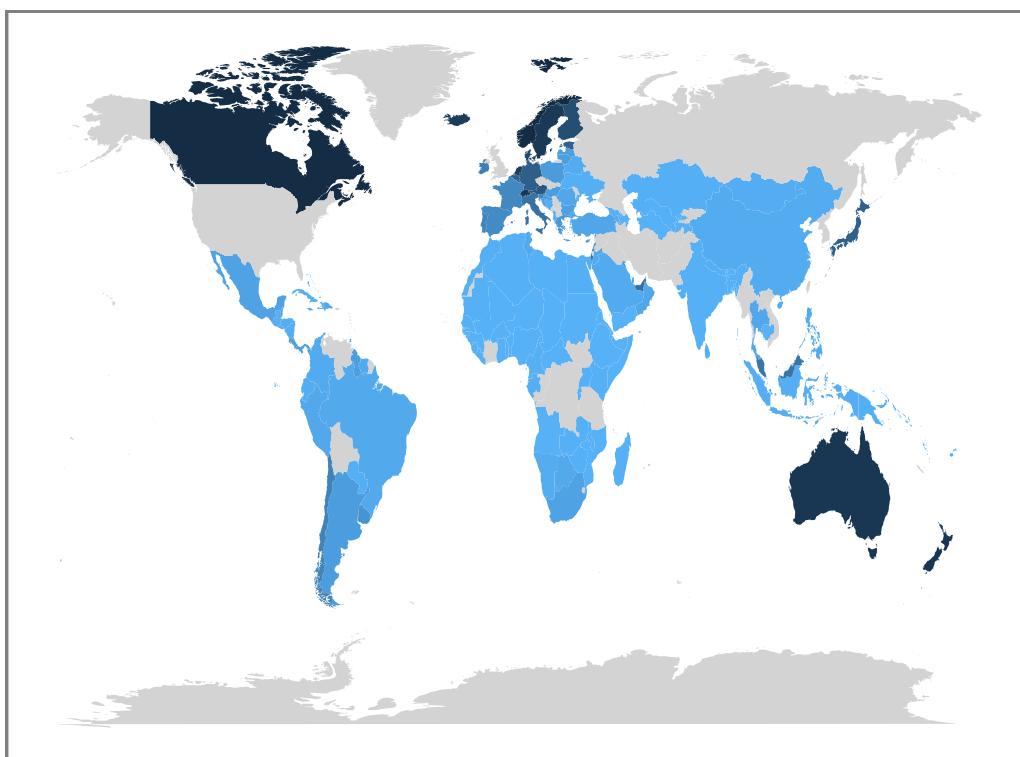
Number of fixed Internet broadband subscriptions, by speed (number)

Year 2000



Internet users per 100 inhabitants

Year 2000



2.3.2 World heatmap - rural vs urban grid

Function to create the world heatmap grid of urban vs rural.

```
### Filter for given indicator, year and loc; delete empty columns
select_ind_year_loc <- function(data, ind, yr, loc) {

  data %>%
    filter(indicator == ind, year %in% yr, location == loc) %>%
    select_if(~sum(!is.na(.)) > 0)

}

### Create world heatmap, rural vs. urban 1x2 grid
facet_heatmap <- function(data, ind, yr) {

  # Create filtered datasets for each location
  data_urban <- select_ind_year_loc(data, ind, yr, "URBAN")
  data_rural <- select_ind_year_loc(data, ind, yr, "RURAL")

  # Joined filtered data with world map
  map_urban <- left_join(map_world, data_urban, by = 'country')
  map_rural <- left_join(map_world, data_rural, by = 'country')

  # Bind all data together
  map <- rbind(map_urban, map_rural)

  # Create data frame with facet property
  map <- data.frame(map, Facet = rep(c("map_urban", "map_rural"),
                                      times=c(nrow(map_urban), nrow(map_rural)))))

  # Text for plot title
  title <- map$seriesdescription %>% na.omit() %>% unique()
  unit <- map$units %>% na.omit() %>% unique() %>% as.character()

  # Plot graph
  ggplot(map, aes(x = long, y = lat, group = group)) +
    geom_polygon(aes(fill = value)) +
    theme(plot.title = element_text(size = 12, face = 'bold'),
          plot.subtitle = element_text(size = 9, hjust = 0),
          axis.ticks = element_blank(),
          axis.text = element_blank(),
          legend.key.size = unit(14, 'pt'))

}
```

```

panel.background = element_rect(fill = 'white'),
panel.border = element_rect(size = 1, colour = 'gray50', fill = NA),
plot.margin = unit(c(0,0,0,11), "mm")) + # top, right, bottom, left

# Faceting
facet_wrap(~Facet)

}

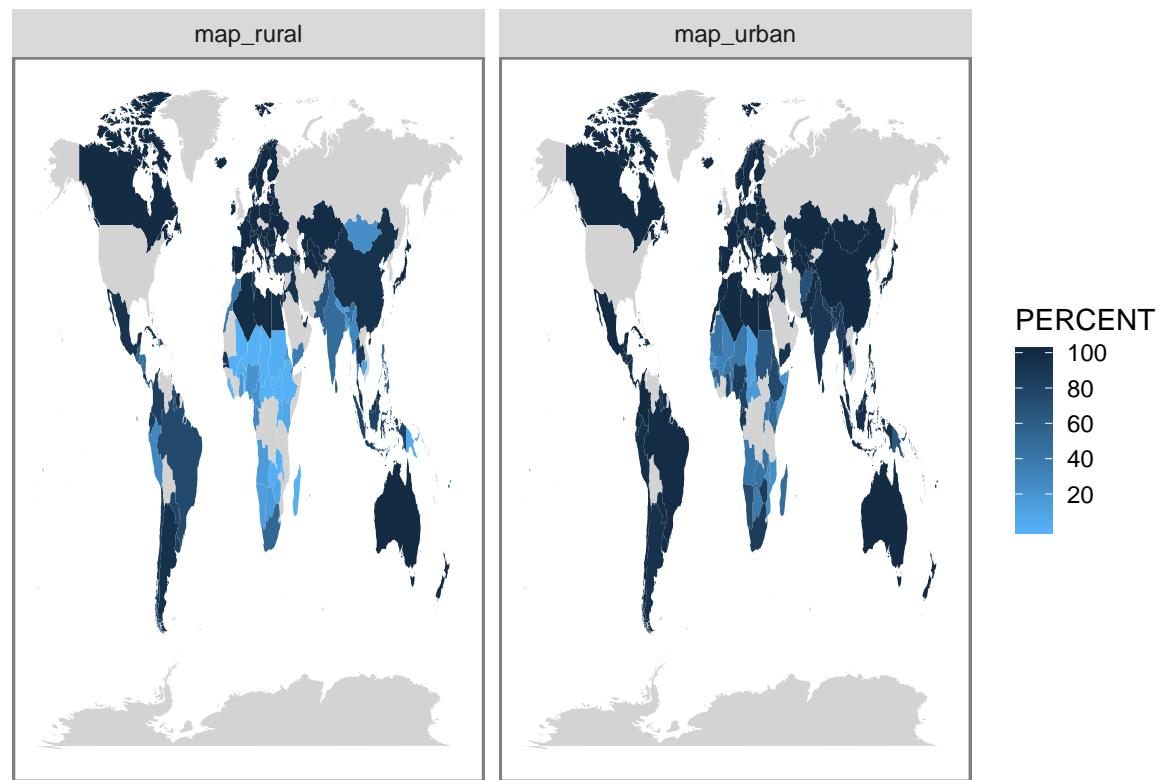
```

Example of world heatmap grid (this plot is more useful in larger dimensions but the concept is still illustrated).

```
facet_heatmap(dat, "7.1.1", 2000)
```

Proportion of population with access to electricity, by urban/rural (%)

Year 2000



2.3.3 Scatterplot - rural vs urban grid

Function for year-to-year scatterplot of the above development, comparing rural to urban. Requires both a start year and an end year as input.

```
### Create scatterplot, rural vs. urban 1x2 grid
facet_scatterplot <- function(data, ind, start, end, x_max = end) {

  # Create function that filters data used for plotting
  filter_data <- function(data, ind, start, end, loc) {
```

```

dat %>%
  # Filter for indicator, given year interval and location
  filter(indicator == ind, year %in% c(start:end), location == loc) %>%
  # Add column for nr. of people covered (in millions) for all countries
  group_by(year) %>%
  mutate(total_people_million = sum(nr_people/10^6, na.rm = TRUE))

}

# Create filtered datasets for each location
rural <- filter_data(data, ind, start, end, "RURAL")
urban <- filter_data(data, ind, start, end, "URBAN")

# Joined filtered data with world map
plot <- rbind(rural, urban)

# Create data frame with facet property
plot <- data.frame(plot, Facet = rep(c("rural","urban"),
                                       times=c(nrow(rural),nrow(urban)))) 

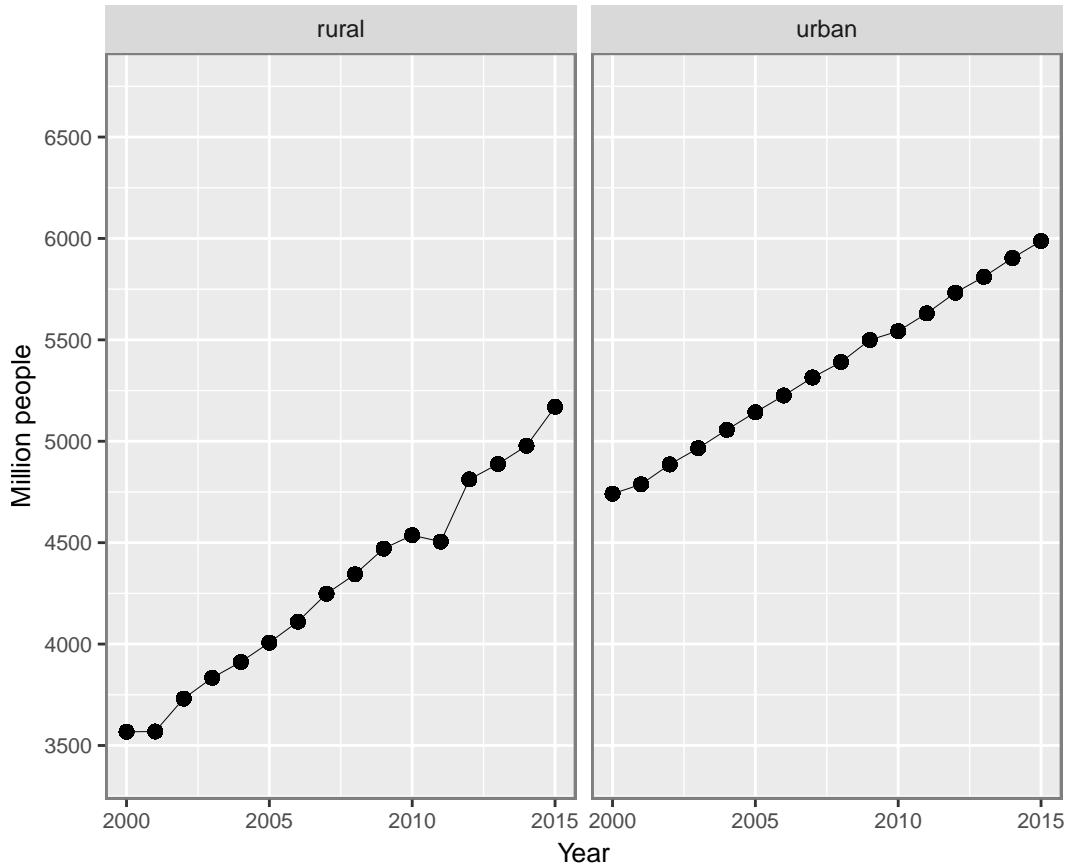
# Specify limits for plot
x_min <- start
#x_max defaults to 'end' but can be overridden via user argument
y_min <- 3400
y_max <- 6750

# Plot graph
plot %>%
  ggplot(aes(year, total_people_million)) +
  geom_point(size = 2) +
  geom_line(size = 0.1) +
  # Set limits
  scale_x_continuous("Year", breaks = seq(2000, 2015, by = 5), limits = c(x_min, x_max)) +
  scale_y_continuous("Million people", breaks = seq(3500, 6500, by = 500), limits = c(y_min, y_max)) +
  # Labels
  labs(x = "Year") +
  # Theme adjustments
  theme(panel.background = element_rect(linetype = 5),
        panel.border = element_rect(size = 1, colour = 'gray50', fill = NA),
        plot.margin = unit(c(0,26,0,0),"mm"), # top, right, bottom, left
        axis.title = element_text(size = 10),
        axis.text = element_text(size = 8)) +
  # Faceting
  facet_grid(~Facet)
}

```

Example of scatterplot.

```
facet_scatterplot(dat, "7.1.1", 2000, 2015)
```



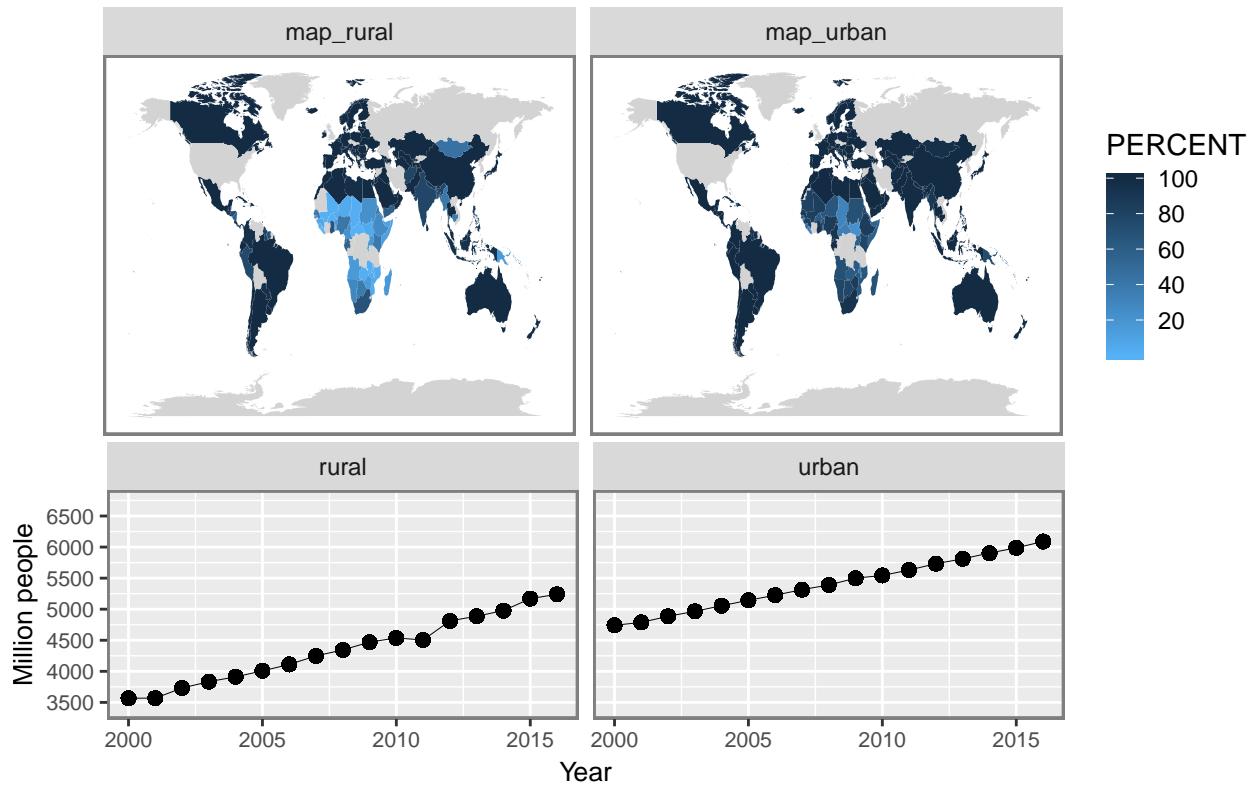
2.3.4 Combined heatmap and scatterplot - rural vs urban grid

It is possible to combine the two grids of plots. Here follows an example for the year 2016 and with the scatterplot starting from year 2000.

```
heatmap <- facet_heatmap(dat, "7.1.1", 2016)
scatterplot <- facet_scatterplot(dat, "7.1.1", 2000, 2016)
grid.arrange(heatmap, scatterplot, nrow = 2, heights = c(0.6, 0.4))
```

Proportion of population with access to electricity, by urban/rural (%)

Year 2016



2.3.5 Animated heatmaps and scatterplots

It is possible to stitch together the above plot to create a time series animation. This is accomplished by stitching the PNGs together into a GIF. An example of the above plot animated in high resolution can be found [here](#). While the functions can be used to create many more animations, a larger selection of examples are found [here](#).

2.3.6 Country bar chart comparison

Below is a function that takes indicator, year and continent as inputs and returns an ordered bar chart of country performance (with the best performing countries first).

```
country_barchart <- function(data, ind, yr, continent) {

  # Filter data according to input arguments and select relevant columns
  filtered_data <- filter(data, year == yr, indicator == ind, region == continent) %>%
    select(c(indicator, seriesdescription, country, year, value, units, location, region))

  # Select only ALLAREA so as to not add up RURAL and URBAN for indicators 6.1.1 and 7.1.1.
  if (ind %in% c("6.1.1", "7.1.1")) {
    filtered_data <- filtered_data %>% filter(location == "ALLAREA")
  }

  # Get bar chart title and unit
```

```

title <- filtered_data$seriesdescription %>% na.omit() %>% unique()
unit <- filtered_data$units %>% na.omit() %>% unique()

# Plot graph
filtered_data %>%
  ggplot(aes(x = reorder(country, -value), y = value)) +
  geom_bar(stat = "identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
  labs(title = paste(title),
       subtitle = paste(continent, '- Year', yr),
       x = "Country",
       y = paste(unit))

}

```

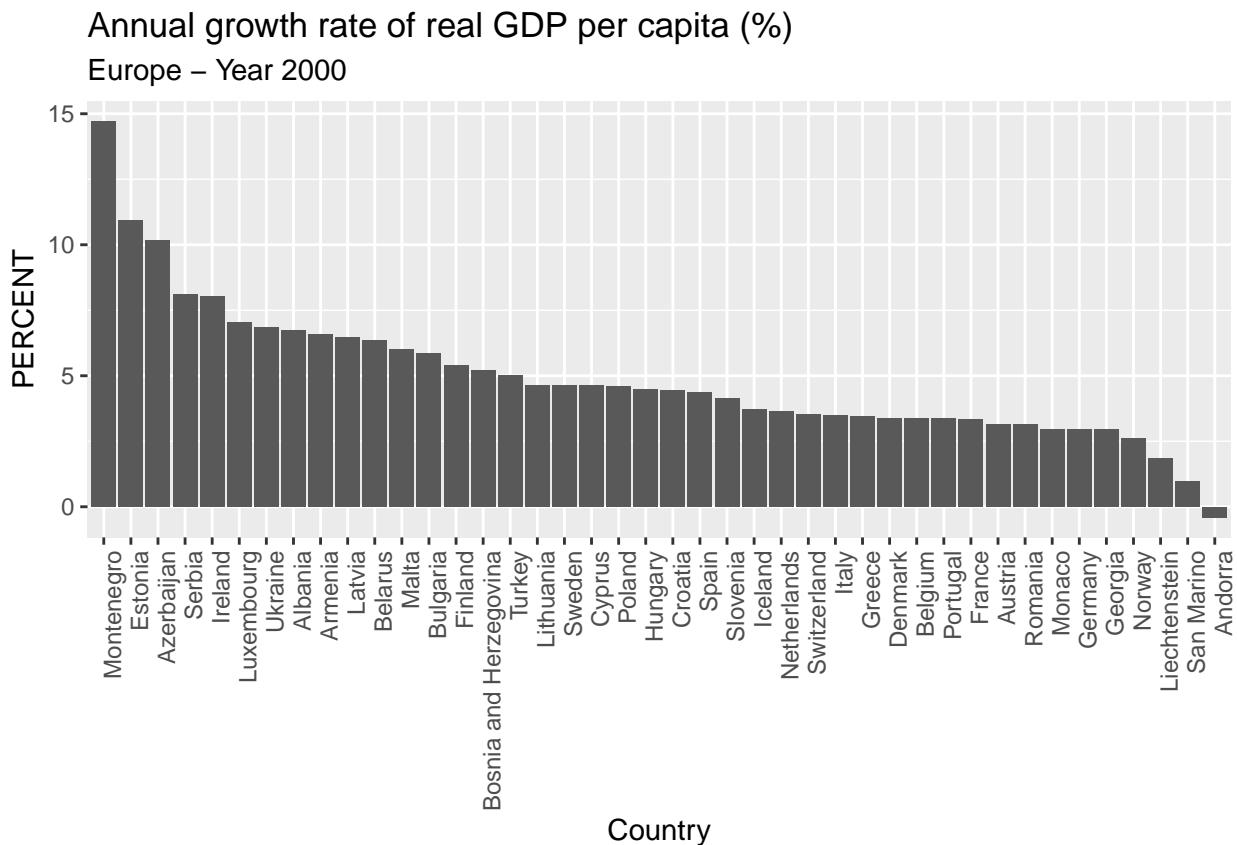
Here are the example charts generated for indicator 8.1.1 (annual RGDP growth) for year 2000, for each continent.

```

continents <- dat$region %>% unique()

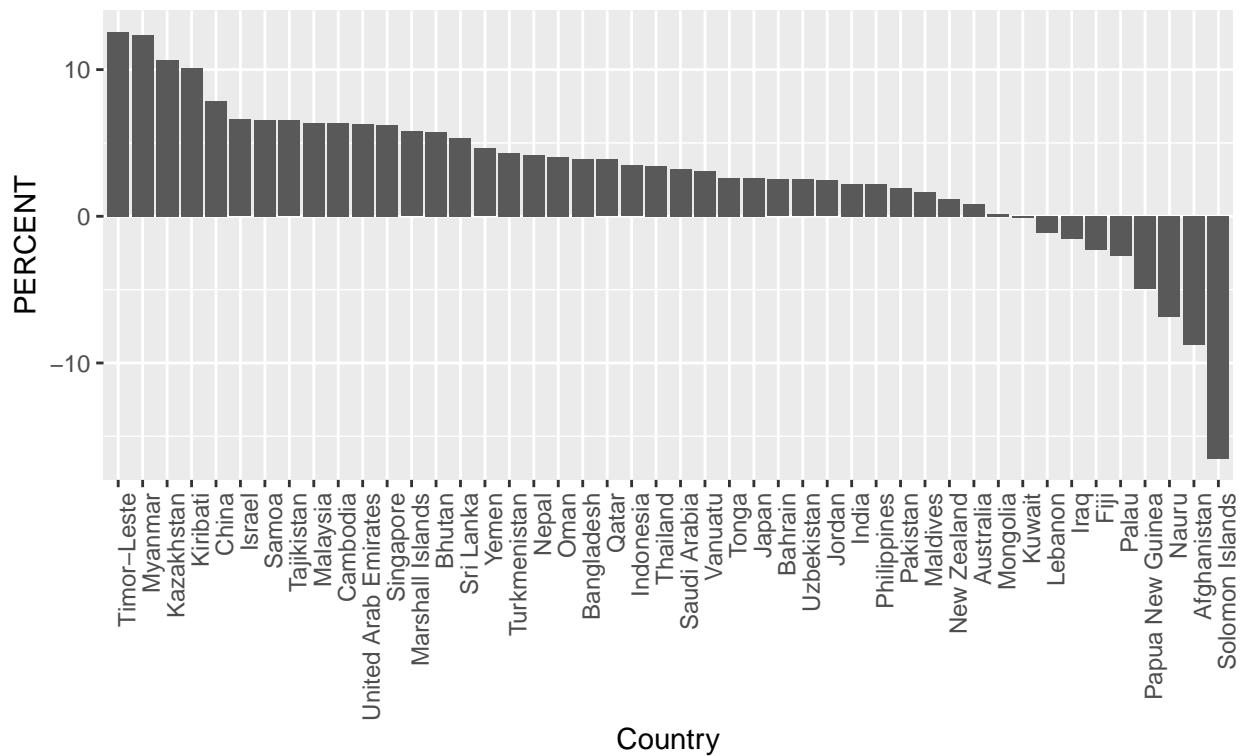
for (cont in continents) {
  plot <- country_barchart(dat, "8.1.1", 2000, cont)
  print(plot)
}

```



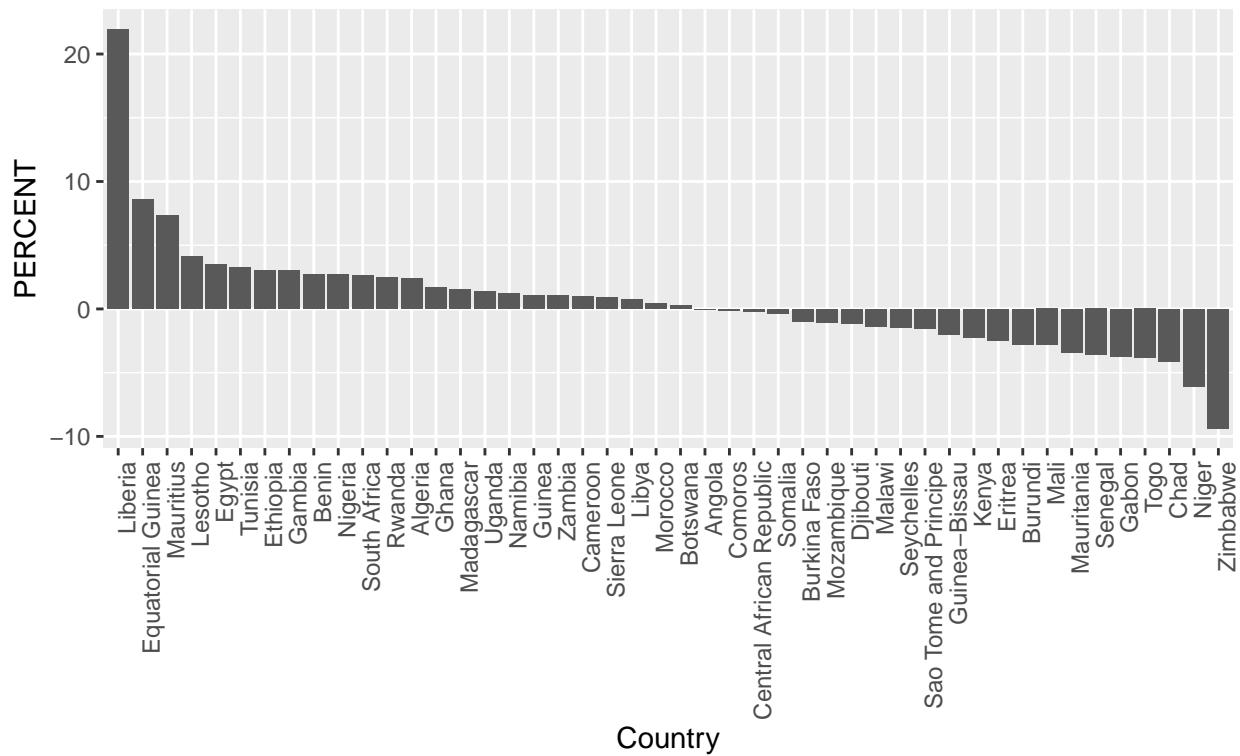
Annual growth rate of real GDP per capita (%)

Asia – Year 2000



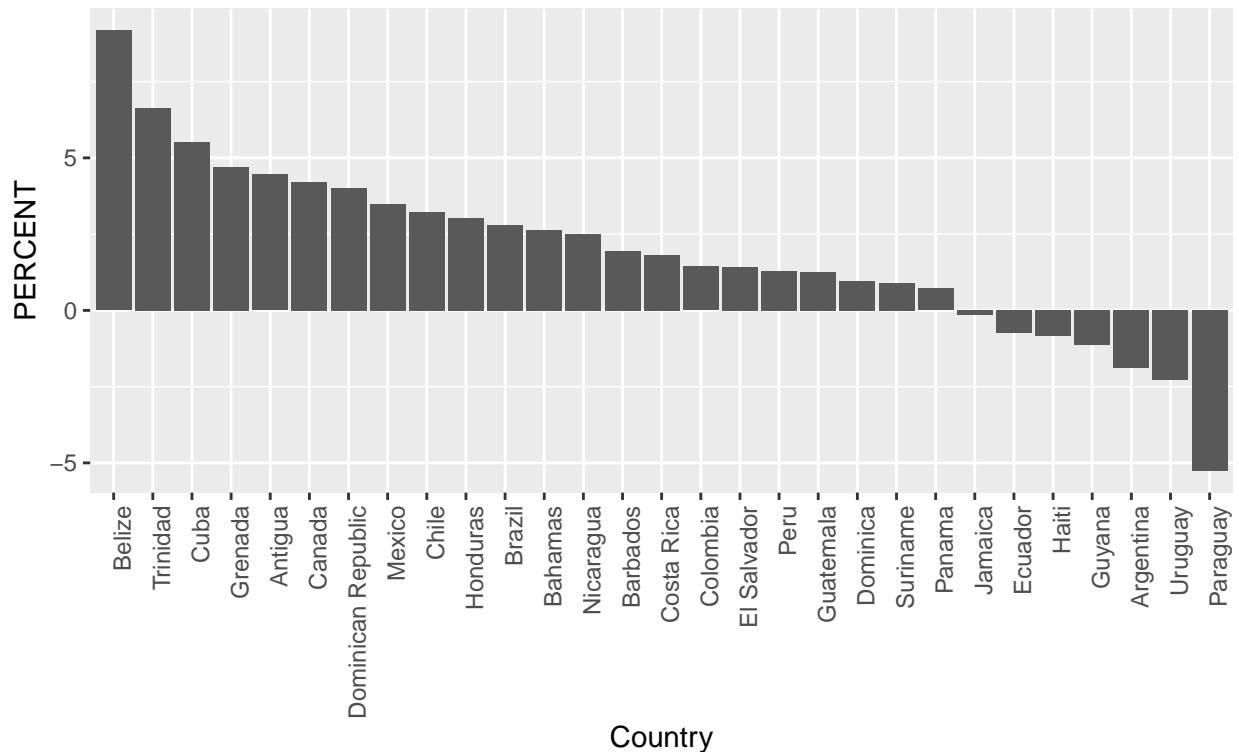
Annual growth rate of real GDP per capita (%)

Africa – Year 2000



Annual growth rate of real GDP per capita (%)

Americas – Year 2000



2.3.7 Wide form time-series table

To inspect indicators numerically, a function was created that takes indicator and continent as input, returning a table in long form for all years there is data.

```
long_form_table <- function(data, ind, continent) {

  # Filter data according to input arguments and select relevant columns
  filtered_data <- data %>%
    select(indicator, seriesdescription, country, year, value, units, location, region) %>%
    filter(indicator == ind, region == continent)

  # Select only ALLAREA so as to not add up RURAL and URBAN for indicators 6.1.1 and 7.1.1.
  if (ind %in% c("6.1.1", "7.1.1")) {
    filtered_data <- filtered_data %>% filter(location == "ALLAREA")
  }

  # Cast data frame to wide form
  dcast(filtered_data, country ~ year)
}
```

Below follows a selection of examples using this function. While all years are recorded in the variable, a smaller selection of years is displayed in the report such that the tables fit on the page.

Proportion of individuals using the Internet in the Americas

2.4.1 Fitting a linear model using lm

Example for Afghanistan on indicator 7.1.1 (population proportion with access to electricity).

```
lm_711_afg <- dat %>%
  filter(country == "Afghanistan", indicator == "7.1.1", location == "ALLAREA") %>%
  select(country, indicator, year, value)

fit <- lm(value ~ year, data = lm_711_afg)

summary(fit)$coeff
```

	Estimate	Std. Error	t value	Pr(> t)
## (Intercept)	-11011.825490	654.756545	-16.81820	3.821371e-11
## year	5.503799	0.326073	16.87904	3.629168e-11

We can use this linear model to predict the future values until 2030, under the assumption of linearity:

```
new <- data.frame(year = 2001:2030)

predict(fit, new)

##          1         2         3         4         5         6
## 1.276348 6.780147 12.283946 17.787745 23.291544 28.795343
##          7         8         9        10        11        12
## 34.299142 39.802941 45.306740 50.810539 56.314338 61.818137
##          13        14        15        16        17        18
## 67.321936 72.825735 78.329534 83.833333 89.337132 94.840931
##          19        20        21        22        23        24
## 100.344730 105.848529 111.352328 116.856127 122.359926 127.863725
##          25        26        27        28        29        30
## 133.367525 138.871324 144.375123 149.878922 155.382721 160.886520
```

The best fit of the model suggests that 100% of Afghanistan's will have access to electricity by 2019. Using the lower confidence bound gives us a more conservative prediction of 2023:

```
predict(fit, new, interval = "prediction")
```

	fit	lwr	upr
## 1	1.276348	-13.966371	16.51907
## 2	6.780147	-8.255178	21.81547
## 3	12.283946	-2.573631	27.14152
## 4	17.787745	3.077195	32.49830
## 5	23.291544	8.696373	37.88672
## 6	28.795343	14.283147	43.30754
## 7	34.299142	19.836959	48.76133
## 8	39.802941	25.357468	54.24841
## 9	45.306740	30.844557	59.76892
## 10	50.810539	36.298343	65.32274
## 11	56.314338	41.719167	70.90951
## 12	61.818137	47.107587	76.52869

```

## 13 67.321936 52.464359 82.17951
## 14 72.825735 57.790410 87.86106
## 15 78.329534 63.086815 93.57225
## 16 83.833333 68.354767 99.31190
## 17 89.337132 73.595543 105.07872
## 18 94.840931 78.810481 110.87138
## 19 100.344730 84.000952 116.68851
## 20 105.848529 89.168334 122.52873
## 21 111.352328 94.313994 128.39066
## 22 116.856127 99.439273 134.27298
## 23 122.359926 104.545469 140.17438
## 24 127.863725 109.633832 146.09362
## 25 133.367525 114.705552 152.02950
## 26 138.871324 119.761759 157.98089
## 27 144.375123 124.803515 163.94673
## 28 149.878922 129.831821 169.92602
## 29 155.382721 134.847611 175.91783
## 30 160.886520 139.851756 181.92128

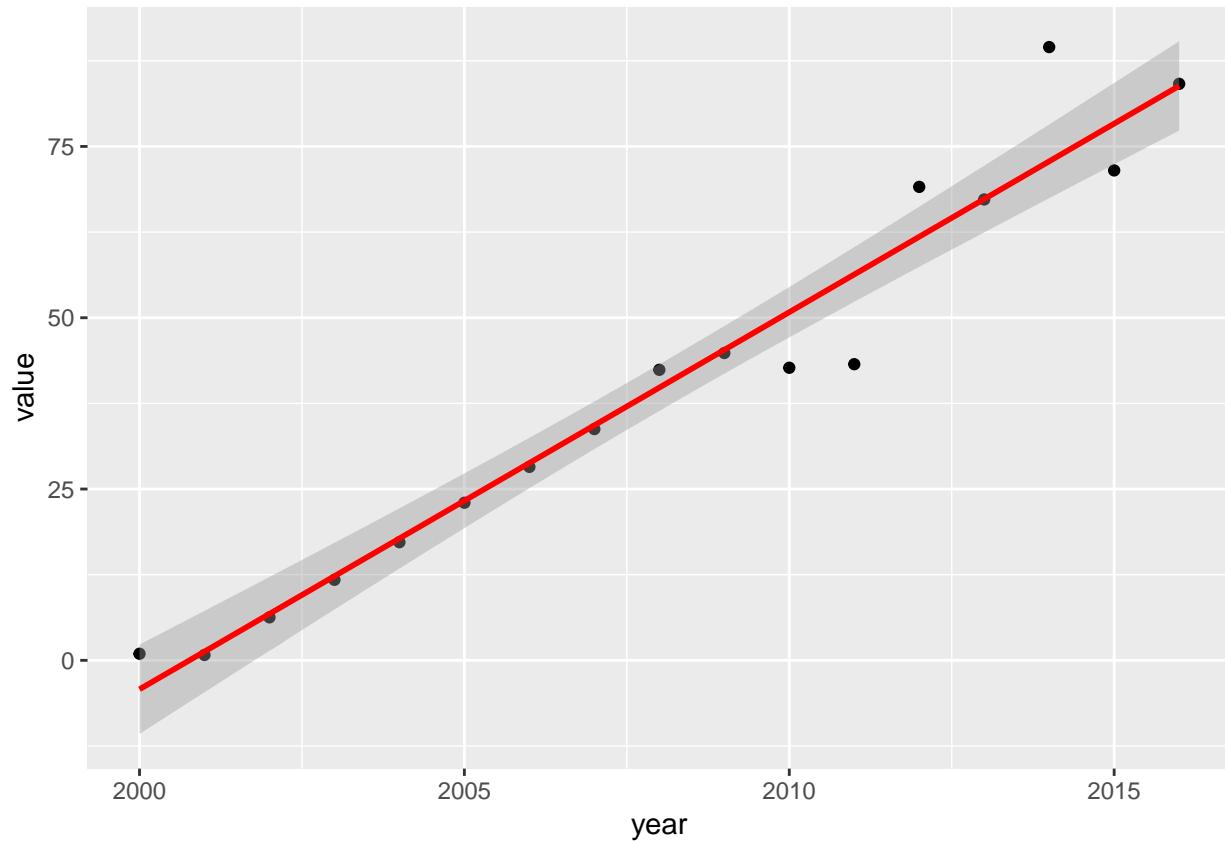
```

A plot of the lm model can be created with ggplot.

```

ggplot(lm_711_afg, aes(x = year, y = value)) +
  geom_point() +
  stat_smooth(method = "lm", col = "red")

```



3 Results

This project has successfully fetched data from all the desired sources and created a unified dataset that is convenient for exploration and analysis. In my opinion, the key achievement of this project is the creation of functions that visualize the results from the dataset. These functions also serve as a baseline for communicating the results of analysis and forecasting. Due to time constraints and the significant effort I have needed to put into the project, I have unfortunately not been able to produce much forecasting. One of the main goals of the projects were to use trend analysis to predict future results. While a simple example of such analysis is included, there is much potential to improve this aspect of the project.

Table 6: Achievement of goals

Goal	Status
Programmatically download data from the relevant sources	Achieved
Clean and merge data into unified datasets	Achieved
Communicate reported results graphically	Achieved ⁽¹⁾
Use trend analysis to predict future results	Lacking

⁽¹⁾ More visualization techniques can always be added.

4 Conclusion

The desired data has been downloaded, cleaned and explored in detail. Visualization of this dataset were considered important from the outset and much effort has been spent on creating good graphics. This has come at the expense of forecasting, which will require more work. The project required more time than anticipated but has been very fun and immensely rewarding; I have learned a lot and only regret not having more time to further the project before submission. I intend to continue working on improving and adding visualizations, and developing forecasting methods including the visualization of those forecasts. I also plan to add a function to pull new data and record the accuracy of the forecasting methods against the new data.