

Towards a programmable Dataverse: Finding, getting and analysing the Data you need

Jan. 2022

Lars Kaczmirek, Lisa Hirsch

Abstract

The presented Open Educational Resources help mitigate common problems that arise when scientific data is used for educational purposes: Often, finding, getting and analyzing data is a time-consuming task. To address this problem, we demonstrate how teachers and their students can use the methods to directly download data. Good analysis scripts begin with code that fetches the data from a repository. This contributes to making research reproducible and supports proper citation of research data. Two examples, one executable in Python, one in a local R installation, will show how to download datafiles from a Dataverse repository, and take first steps in analysing data.



To the extent possible under law, Lars Kaczmirek and Lisa Hirsch have waived all copyright and related or neighboring rights to Towards a programmable Dataverse: Finding, getting and analysing the Data you need. this work is published from: Austria.

When using these resources we kindly ask you to cite our work: Kaczmirek, L. & Hirsch, L. (2022). Towards a programmable Dataverse: Finding, getting and analysing the data you need. Vienna: AUSSDA. Downloaded from: <https://github.com/AUSSDA/dataverse-oer-resources>

This CC0 waiver does not cover the research data or software packages that are being used and referenced in the code. Trademark and patent rights are not affected. Dataverse, Python, R, SPSS, Stata are registered trademarks and we are not endorsed by the rights holders. Please make sure to also appropriately cite all research data that you use.

1 Introduction

AUSSDA - The Austrian Social Science Data Archive is a rich source of data for research purposes. Some data can also be used for other purposes under the specified conditions. In these materials which AUSSDA distributes freely we show how to access the AUSSDA Dataverse via different programming tools, and download data that is open access. We also show how to dive into data analysis by introducing summary statistics and plotting data.

Imagine the following ideal situation: A teacher in a master class discusses a scientific article with students. They do so by re-running and making changes to the analysis code that was distributed with the published article. The programming code automatically fetches the data from the repository that was used by the authors of the article to archive and share the respective data. Students are enabled to reproduce the results including tables and figures. As the teacher makes suggestions to change the code students are encouraged to increase their analysis and critical thinking skills.

Compare the previous scenario to the often experienced situation in which teachers or students have to register to gain access to data and time is spent on delivering the data to the students. Often more time is lost in figuring out how downloaded data has to be prepared before the above described analyses code runs without errors. Thus the current situation often means that teachers have to spend extra hours to prepare article, analysis or replication code and the corresponding data before the work with the students can begin.

These materials help mitigate the problems that arise from the scenario above: As a teacher, you use the direct data download methods in your teaching to make sure that everyone uses the correct files. No more wasted time with distributing data files to a class and checking on files not found. As a student, starting your analysis script with code to download the data goes a long way to make your research reproducible. No more confusion about which data set or version has actually been used to run the analysis.

The materials are hosted on <https://github.com/AUSSDA/dataverse-oer-resources>. It contain the public domain dedication under which this work is published, a general readme document, this README.pdf, a Jupyter Notebook to run the example in Python via mybinder.org and a file that passes the requirements to mybinder.org, an R script to execute the example using a local R installation and a pdf that explains the R code.

Whenever you reuse data from others you have the moral obligation and in many cases also the legal obligation to cite your sources appropriately. Remember to always include the persistent identifier - for data in the AUSSDA Dataverse this is the DOI. You can also refer to <https://dataverse.org/best-practices/data-citation> for best practices.

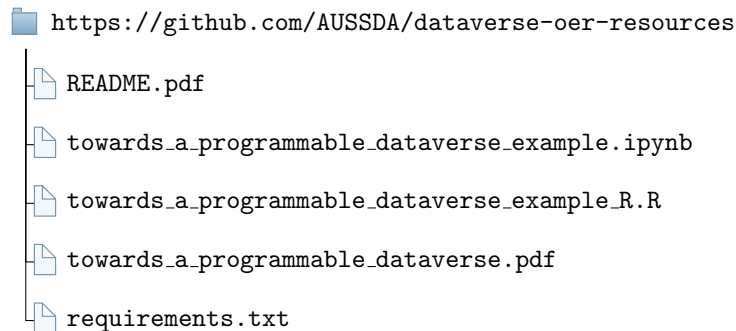


Figure 1: Structure of the materials hosted on Github

2 Using the Dataverse with Python

You will learn how to navigate to the repository and download a datafile, and how to use different datafiles.

```
1 # This installs the packages so Jupyter Notebook can execute the
  code.
2 # Version info: python 3.9.7
3 # If you run this in a mybinder.org environment, you do not need to
  execute this block of code.
4 import sys
5 !conda install --yes --prefix {sys.prefix} requests
6 !conda install --yes --prefix {sys.prefix} jsonschema
7 !conda install --yes --prefix {sys.prefix} pandas
8 !conda install --yes --prefix {sys.prefix} numpy
9 !conda install --yes --prefix {sys.prefix} matplotlib
10
11 ![sys.executable] -m pip install pyDataverse
```

Start by installing and loading the packages you need.

```
1 # The package pyDataverse handles everything you need to connect
  and download files from a Dataverse installation.
2 # Check out the documentation:
3 # https://pydataverse.readthedocs.io/en/latest/user/basic-usage.
  html#download-and-save-a-dataset-to-disk
4
5 from pyDataverse.api import NativeApi, DataAccessApi
6 import io
7 import pandas as pd
```

Decide which data repository you want to use, and which dataset you want to download. This example is based on the AUSSDA Dataverse. AUSSDA - The Austrian Social Science Data Archive is a rich source for data for research purposes. Some data can also be used for other purposes under the specified licence conditions.

You will learn how to download the Open Access edition of the Social Survey

Austria 2016 (Bacher et al., 2018). First you define the repository URL, then the persistent identifier (DOI).

```
1 # Step 1: Define and connect to repository API
2 repository_url = 'https://data.aussda.at/'
3 native_api = NativeApi(repository_url)
4
5 # Step 2: Define and download dataset
6 DOI = "doi:10.11587/EHJHFJ"
7 dataset = native_api.get_dataset(DOI)
```

To download the data file automatically, the following code proceeds in several steps. First, you create a list of available files. Then, you download the first datafile. In Dataverse, there are several files stored under each DOI: datafiles in various formats suitable for STATA, SPSS, R, tab-delimited, and documentation as pdf-files. The code first goes through all files in the file list (i.e. all files in the dataset), and identifies all files that are tab-delimited files. These are the datafiles you want to download when using Python. Then, it picks the first of these tab-delimited files which is usually the original datafile by default. Then you download this file.

```
1 # Create a list of files that are available in the dataset
2 files_list = dataset.json()['data']['latestVersion']['files']
3
4 # Find the first tab-file
5 ident = []
6
7 for file in files_list:
8     filename = file["dataFile"]["filename"]
9     file_id = file["dataFile"]["id"]
10    if filename.endswith('.tab'):
11        ident.append(file_id)
12
13 # Save the ID of the first tab file that occurs in the list
14 datafile_id = ident[0]
15
16 # Step 3: Connect to API for data access at the repository
17 data_access_api = DataAccessApi(repository_url)
18
19 # Step 4: Download data file using its id
20 response = data_access_api.get_datafile(datafile_id)
```

Now that you have downloaded the datafile, you can work with it: It needs to be transformed into a format that can be used by the Pandas package.

```
1 ## Working with the datafile Pt. 1
2
3 # Transform the response into a Pandas data frame
4 data = io.StringIO(str(response.content, 'utf-8'))
5
6 # The file is a tab-delimited file in the repository, i.e. the
   separator between columns (variables) is a tab.
7 # The first line (the header) consists of the variable names.
8 data = pd.read_csv(data, sep="\t", index_col=0)
```

In the next step you check if the download worked, and what the data looks like.

The rows are unique observations. In your case each row is an individual respondent, because you use survey data. The columns are variables, which are mostly questions in survey data. The variable names in the header show which question is stored in which column, and the values show the answer each respondent gave to each question.

```
1 ## Working with the datafile Pt. 2
2
3 # Show sample data by displaying the first rows
4 data.head()
```

Now that you know how to download a dataset you may want to try a different dataset that is also published under an open access license such as CC BY that does not require a login or API key.

```
1 # If you want to use the survey data of the Corona Panel Project (
  Kittel et al, 2020) change the DOI in the code above to
  "10.11587/P5YJ00" in Step 2. Look at the codebook or the
  questionnaire to determine which variable contains information
  on which question.
```

There is also a way to access a file directly: If the automatic download described above does not leave you with the desired file, you can specify the name of the file you want to retrieve manually, and download it using the function `get_dataframe_by_name`. In order to find the name, check the object `files_list` that you created after Step 2, or look at the landing page for the Social Survey Austria in the repository at <https://doi.org/10.11587/EHJHFJ>.

```
1 # Manual Process: Specify the id of the file you want to download,
  and proceed with Step 1 and Step 2
2 # Check the id numbers for each file.
3 for file in files_list:
4     filename = file["dataFile"]["filename"]
5     file_id = file["dataFile"]["id"]
6     print("File name {} has id {}".format(filename, file_id))
7
8 # Specify which file you want to retrieve.
9 datafile_id = 188
10 # Step 1: Connect to API for data access at the repository
11 data_access_api = DataAccessApi(repository_url)
12 # Step 2: Retrieve data file
13 response = data_access_api.get_datafile(datafile_id)
```

3 Using the Dataverse with R

You will learn how to navigate to the repository and download a datafile, and how to use different datafiles.

In order to execute the code, you need an installation of the software R, for example RStudio. Please refer to <https://www.rstudio.com/products/rstudio/download/> for more information on how to install RStudio. You can then copy the code below into a R Script, and execute the commands locally.

Start by installing and loading the packages you need. If you have not already installed the packages, remove the comment signifier in the first line to execute the function `install.packages`.

```
1 # install.packages("dataverse", "tibble", "ggplot2")
2
3 library(dataverse)
4 library(tibble)
5 library(ggplot2)
```

Decide which data repository you want to use, and which dataset you want to download. This example is based on the AUSSDA Dataverse.

You will learn how to download the Social Survey Austria 2016 (Bacher et al., 2018) from the Dataverse which is open access. First you define the repository URL, then the persistent identifier to the dataset (DOI).

```
1 # Step 1: Define and connect to repository API
2 repository_url <- "https://data.aussda.at/"
3 Sys.setenv("DATAVERSE_SERVER" = repository_url)
4
5 # Step 2: Define dataset and retrieve dataset
6 DOI <- "doi:10.11587/EHJHFJ"
7 dataset <- get_dataset(DOI)
```

To download the data file automatically, the following code proceeds in several steps: First, you create a list of available files. Then, you download the first datafile.

In Dataverse, there are several files saved under each DOI: datafiles (usually .tab files), and documentation as pdf-files. The following code saves a list of all these files.

```
1 # Downloading Data: list files
2 files_list <- dataset$files[c("filename", "contentType")]
```

The following code first goes through all files in the file list (i.e. all files in the dataset), and identifies all files that are tab-delimited files. In Dataverse, there are several files stored under each DOI: datafiles in various formats suitable for STATA, SPSS, R, tab-delimited, and documentation as pdf-files. The code first goes through all files in the file list (i.e. all files in the dataset), and identifies all files that are tab-delimited files. These are the datafiles you want to download when using R. Then you download this file.

```
1 # Downloading Data: identify desired file
2 filename <- files_list$filename
3 ext <- unlist(strsplit(filename, "\\.")) [c(F,T)]
4 ident <- filename[which(ext %in% "tab")][1] # save the filename of
      the first .tab file
```

Use the function `get_dataframe_by_name` from the package `Dataverse` to download the datafile directly.

```
1 # Downloading Data: download data
2 data <- get_dataframe_by_name(
3   paste(ident),
4   paste(DOI),
5 )
```

The function `get_dataframe_by_name` already reads in the data in the format that can be used by R.

The rows are unique observations. In your case each row is an individual respondent, because you use survey data. The columns are variables, which are mostly questions in survey data. The variable names in the header show which question is stored in which column, and the values show the answer each respondent gave to each question.

```
1 # Show sample data by displaying the first rows
2 head(data)
```

After having downloaded data, you can now start to analyse the data. In order to know what each variable means, and what its values mean, look at the questionnaire or into the codebook for the dataset you are retrieving. Go to <https://doi.org/10.11587/EHJHFJ> and look for a pdf file called codebook. There you will find information on each variable, for example the question that was asked, and possible answers.

The code below is an example of how to conduct some data analysis. The variable `bf29_7` holds information on how much importance respondents assign to the societal importance of their job in their professional lives. (german wording: "Geben Sie auf der folgenden Liste bitte jeweils an, für wie wichtig Sie persönlich das für die berufliche Arbeit halten. Ein Beruf, der für die Gesellschaft nützlich ist."). The answers range from "1 = very important" to "5 = not important at all".

First, in order to exclude respondents who did not or could not answer the question, you set their responses to missing.

```
1 # Plotting Data
2 # recode missings and "cannot say" as missing values
3 # -99 = respondent has not worked before
4 # 8 = cannot say
5 data$bf29_7[data$bf29_7==-99] <- NA
6 data$bf29_7[data$bf29_7==8] <- NA
```

The code below shows how to produce a boxplot that shows the distribution of replies in relation to the respondents' gender.

```
1 # Where to go from here? Pt. 1: Plotting Data
2 boxplot(bf29_7~SEX_2016,data=data,
3         horizontal=TRUE,
4         names=c("male","female"),
5         col=c("wheat","thistle"),
6         xlab="1=very important, 5=not important at all
7         (without cannot say)",
8         ylab="Gender",
9         main="Importance of societal relevance of one's job")
```

The code below shows how to produce a grouped bar plot of the replies in relation to the respondents' gender (see Figure 2).

```

1 # Plotting Data
2 # ggplot for bar plot by groups
3 # define facet labels
4 to_string <- as_labeller(c('1' = "male", '2' = "female"))
5
6 # plot
7 myplot <- ggplot(data=data, aes(x= bf29_7, group=SEX_2016)) +
8   geom_bar(aes(y = ..prop.., fill = factor(..x..), stat="
9     count") +
10     scale_y_continuous(labels=scales::percent) +
11     labs(y = "Percent",
12         x="1=very important, 5=not important at all, without
13         cannot say",
14         title="Importance of societal relevance of one's job",
15         fill="Wichtigkeit") +
16     guides(fill=FALSE)+
17     facet_grid(~SEX_2016, labeller = to_string)
18 myplot

```

Another option is to use a different dataset that is published under an open access license (CC BY) that does not require a login or API key.

If you want to use the Corona Panel Project (Kittel et al., 2020) instead, change the DOI in the code above to "10.11587/P5YJ0O" in Step 2. Look at the codebook or the questionnaire again to determine which variable contains information on which question.

If the automatic download described above does not work, you can specify the name of the file you want to retrieve manually, and download it using the function `get_dataframe_by_name`. In order to find the name, check the object `files_list` that you created after Step 2, or look at the landing page for the Social Survey Austria in the repository at <https://doi.org/10.11587/EHJHFJ>.

You can also download the files manually when using the Corona Panel Project in case the automatic download does not work.

```

1 # Downloading data manually by accessing a file directly
2 data <- get_dataframe_by_name(
3   "10007_da_de_v1_2-1.tab",
4   paste(DOI))

```

4 Using the Dataverse with SPSS and STATA

You can also download files from the Dataverse using other statistical software packages like SPSS and STATA.

References

Bacher, J., Beham-Rabanser, M., Grausgruber, A., Haller, M., Höllinger, F., Muckenhuber, J., ... Verwiebe, R. (2018). *Social Survey Austria 2016*.

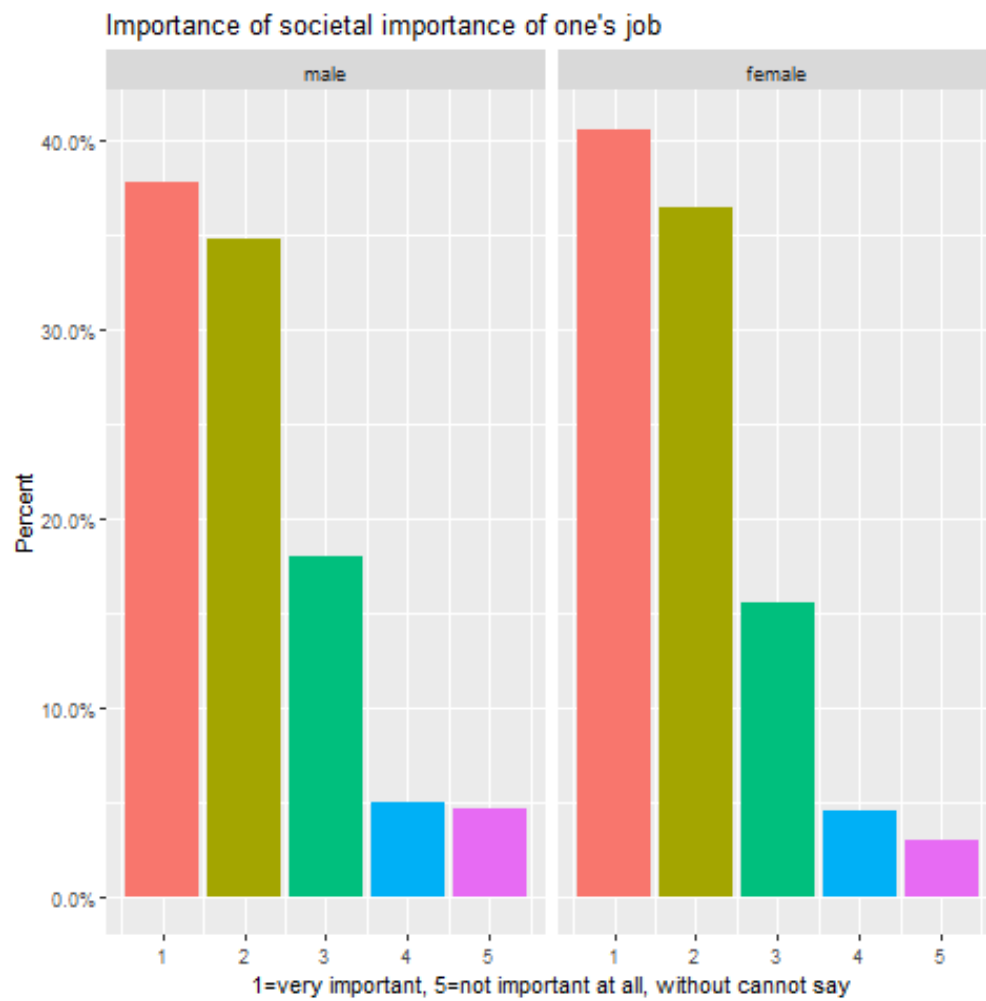


Figure 2: Example: Distribution of societal relevance of one's job

<http://dx.doi.org/10.11587/EHJHFJ>. (V3, AUSSDA) doi: 10.11587/EHJHFJ

Kittel, B., Kritzinger, S., Boomgaarden, H., Prainsack, B., Eberl, J.-M., Kalleitner, F., ... Schlogl, L. (2020). *Austrian Corona Panel Project (OA edition)*. <http://dx.doi.org/10.11587/P5YJ00>. (V2, AUSSDA) doi: 10.11587/P5YJ00