1. File formats

The Enroll-HD PDS dataset is provided in two formats:

- **CSV file:** CSV stands for comma separated values (.csv) which is a delimiter-separated format. The PDS data uses the **tab** as the delimiter. Software settings need to be adapted respectively.

- **R file**: binary code format for the R[1] software application (a software environment for statistical analysis).

Because of the complexity and the size of the data set, use of a statistical software package such as R, Stata, or SAS is recommended. The .csv file format can also be imported into Excel (caution is advisable).

It is important that files **are not be edited in a word processing software or other programs that may potentially modify characters**, as this may damage the integrity of the original files. CSV files can be saved in other formats which are compatible with other statistical software packages as needed.

---

[1] R Core Team (2017). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

2. Importing data

### *Importing CSV files into Excel*

The .csv files can be imported and opened in Microsoft Excel. Because Excel is language dependent and delimiters differ from one country to another, some considerations need to be addressed when opening the .csv files to maintain data integrity. The procedures outlined here, to open the .csv files, can be applied to most recent versions of Excel.

As a default, Excel reads the values for each column as being in a "**General**" format. For example, unless otherwise specified, Excel interprets numeric data as numbers (e.g., 1234), entered dates as date format (as pre-set, e.g. 11/28/2016), and changes other values (e.g. strings) to text format (e.g. Aspirin). For some entries this is counterproductive, as Excel may **misinterpret entries** and **incorrectly reformat the data**, effectively changing the data (e.g. 1.5 is read as May 1 instead of 1.5 mg; or the WHO-DD Code for Tetrabenazine 00222101003 is changed to 22211003, removing the important leading "0"s).

To maintain the integrity of the data, each data column needs to be **carefully examined prior to importing the data into Excel**.

**An illustrated guide** for correctly importing CSV data files into Excel are provided in **Appendix A.**

### *Importing CSV files into R*

Make sure the CSV file has not been opened and saved using a word processing software. A software package capable of reading CSV files must be loaded into R environment. The package "readr" is one of the most popular packages, but there are several others that will also work. If a package like "readr" is not already installed, the CSV data files can be imported using the following code line:

```
install.packages(readr)
```

To load the CSV data into R using a package like "readr" use the: **library(readr)** command. To ensure the CSV file is imported correctly, set the directory to the file folder where the PDS files are located, and then run the following code:

```
file = read_delim("file.csv", "\t", escape_double = FALSE, trim_ws = TRUE)
```

### Importing R files into R

This data file is specific for R. After loading the R data files into R, 9 data frames are made available in the R environment and are ready to be used. The loading can be done using the function command:

```
load("Rdata_directory")
```

For **Rstudio** users, the loading can be performed by clicking in the "load workspace" ribbon, and then browsing for the location of the R data file.

## Appendix A: An illustrated guide to correctly importing CSV files into Excel

The file used for this demonstration is the 'profile.csv' file.

**Step 1 – Open CSV file in Excel:** Open the .csv file using Excel, or open Excel and on the "Data" tab click "**From Text/CSV**". Data will be imported in entirety into the first column of the Excel file, as illustrated below.

**Step 2 – Open Text to Columns Wizard:** Select the first column, then on the tab "Data" click "**Text to Columns**". A wizard will appear to guide you through the process.

**Step 3 – Specify data file type:** In the Text to Columns Wizard (step 1 of 3), select the "**Delimited**" checkbox (this lets Excel know that the data fields are separated by commas or tabs), then click "**Next**".

**Step 4 – Select delimiter type:** In the Text to Columns Wizard (step 2 of 3), select the Delimiter type "**Tab**" (this lets Excel know that the data fields are separated by tabs specifically), then click "**Next**".

**Step 5 – Assign column formats:** For each column (i.e., variable), an appropriate format needs to be assigned. This is completed in the Text to Columns Wizard (step 3 of 3). The default format "**General**" works for most columns. Columns where numbers have leading "0" and columns with mixed entries like 1.5, 1,5, 1/5, need to be explicitly formatted as "**Text**", as entries might otherwise become corrupted in an unchangeable way. After assigning the correct format to each column, click "**Finish**".



**NB:** The data files *pharmacotx* and *nutsuppl* contain two columns '*cmtrt_decod*' and '*cmdostot*' that require formatting as "**Text**".

**Step 6 – Save data file:** The .csv file is now column-separated and should be saved as an Excel file (.xls or .xlsx) using the 'Save As' option.

| | subjid | region | sex | race | handed | hxsid | dssage | dsplace | dsend | caghigh | caglow | momhd | momagesx | dadhd | dadagesx | fhx | ccmtr | ccmtrage | sxsubj | sxfam | hddiagn | sxest | sxrater |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | R0000245 | Northern A | m | 1 | 1 | 0 | | | | 44 | 19 | 1 | 50 | 0 | | 1 | 1 | 42 | 42 | 42 | 47 | 1 | 42 |
| 3 | R0002172 | Europe | f | 1 | 1 | 0 | | | | 38 | 24 | 1 | 55 | 0 | | 1 | 0 | NA | NA | | | | |
| 4 | R0002394 | Europe | m | 1 | 1 | 0 | | | | 41 | 20 | 0 | | 1 | 70 | 1 | 1 | 65 | 67 | 70 | 70 | 1 | 66 |
| 5 | R0007988 | Europe | f | 1 | 1 | 0 | | | | 44 | 20 | 0 | | 1 | 30 | 1 | 0 | NA | NA | | NA | | |
| 6 | R0010845 | Northern A | f | 1 | 1 | | | | | 22 | 15 | | | | | 1 | | | | | | | |
| 7 | R0011592 | Europe | f | 1 | 1 | 0 | | | | 17 | 16 | 0 | | 1 | 40 | 1 | 0 | | | | | | |
| 8 | R0012253 | Europe | m | 1 | 1 | 0 | | | | 41 | 17 | 0 | | 1 | 38 | 0 | 1 | 45 | 45 | 44 | 46 | 1 | 44 |
| 9 | R0012861 | Northern A | f | 6 | 1 | 0 | | | | 46 | 17 | 1 | 32 | 0 | | 0 | 0 | NA | NA | | NA | | |
| 10 | R0017074 | Europe | m | 1 | 1 | 0 | | | | 47 | 18 | 1 | 25 | 0 | | 1 | 1 | 35 | 40 | 35 | 40 | 1 | 35 |
| 11 | R0017175 | Northern A | m | 1 | 1 | 1 | | | | 27 | 17 | 1 | 56 | 0 | | 1 | 0 | | | | | | |
| 12 | R0018320 | Europe | m | 1 | 1 | 1 | | | | 42 | 18 | 1 | 35 | 0 | | 0 | 1 | 65 | 65 | 65 | 68 | 1 | 65 |
| 13 | R0019242 | Europe | | 8 | 1 | | | | | 23 | 16 | | | | | 0 | | | | | | | |
| 14 | R0022010 | Northern A | f | 1 | 2 | 0 | | | | 44 | 14 | 1 | 40 | 0 | | 1 | 1 | 49 | 49 | 49 | 52 | 1 | 49 |
| 15 | R0022157 | Northern A | f | 1 | 1 | 0 | | | | 18 | 17 | 1 | 50 | 0 | | 0 | 0 | NA | | | | | |
| 16 | R0023119 | Northern A | f | 15 | 1 | 0 | | | | 40 | 16 | 0 | | 0 | | 1 | 1 | 61 | | | | | |
| 17 | R0023306 | Europe | m | 1 | 1 | 0 | | | | 41 | 16 | 1 | 50 | 0 | | 1 | 0 | | | | | | |
| 18 | R0025704 | Europe | m | 1 | 1 | | | | | 20 | 16 | | | | | 1 | | | | | | | |
| 19 | R0025936 | Europe | m | 1 | 1 | 0 | | | | 42 | 20 | 1 | 60 | 0 | | 1 | 0 | NA | NA | | | | |
| 20 | R0027760 | Europe | f | 1 | 1 | | | | | 23 | 17 | | | | | 1 | | | | | | | |
| 21 | R0027983 | Northern A | f | 1 | 1 | 0 | | | | 42 | 21 | 0 | | 1 | 40 | 1 | 1 | 54 | 54 | NA | | NA | |
| 22 | R0028347 | Europe | m | 1 | 1 | 0 | | | | 43 | 20 | 0 | | 1 | 50 | 1 | 1 | 51 | 48 | 48 | 52 | 1 | 48 |
| 23 | R0029018 | Northern A | f | 1 | 1 | 0 | | | | 20 | 19 | 0 | | 1 | 43 | 1 | 0 | | | | | | |
| 24 | R0029189 | Northern A | f | 2 | 1 | 1 | | | | 16 | 16 | 0 | | 1 | 44 | 1 | 0 | NA | NA | | | 0 | |
| 25 | R0029959 | Northern A | m | 1 | 1 | | | | | 18 | 18 | | | | | 1 | | | | | | | |
| 26 | R0029995 | Northern A | f | 1 | 1 | 1 | | | | 45 | 17 | 1 | 40 | 0 | | 1 | 0 | | | | | | |
| 27 | R0030050 | Europe | f | 1 | 1 | 1 | | | | 42 | 16 | 1 | 42 | 0 | | 1 | 1 | 44 | 44 | NA | 45 | 1 | 45 |
| 28 | R0030270 | Northern A | f | 1 | 1 | 0 | | | | 43 | 26 | 1 | 49 | 0 | | 1 | 0 | NA | NA | | NA | | |
| 29 | R0031143 | Europe | f | 1 | 1 | 0 | | | | 45 | 16 | 0 | | 1 | 53 | 1 | 1 | 35 | 35 | 35 | 38 | 1 | 35 |
| 30 | R0031614 | Europe | f | 1 | 1 | 0 | | | | 45 | 19 | 0 | | 1 | 40 | 1 | 1 | 31 | | | 31 | 1 | 31 |
| 31 | R0031839 | Northern A | f | 1 | 2 | 0 | | | | 37 | 18 | 1 | 55 | 0 | | 1 | 1 | 54 | 54 | 54 | | 0 | |