

## 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<sup>1</sup> 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.

---

<sup>1</sup> 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:

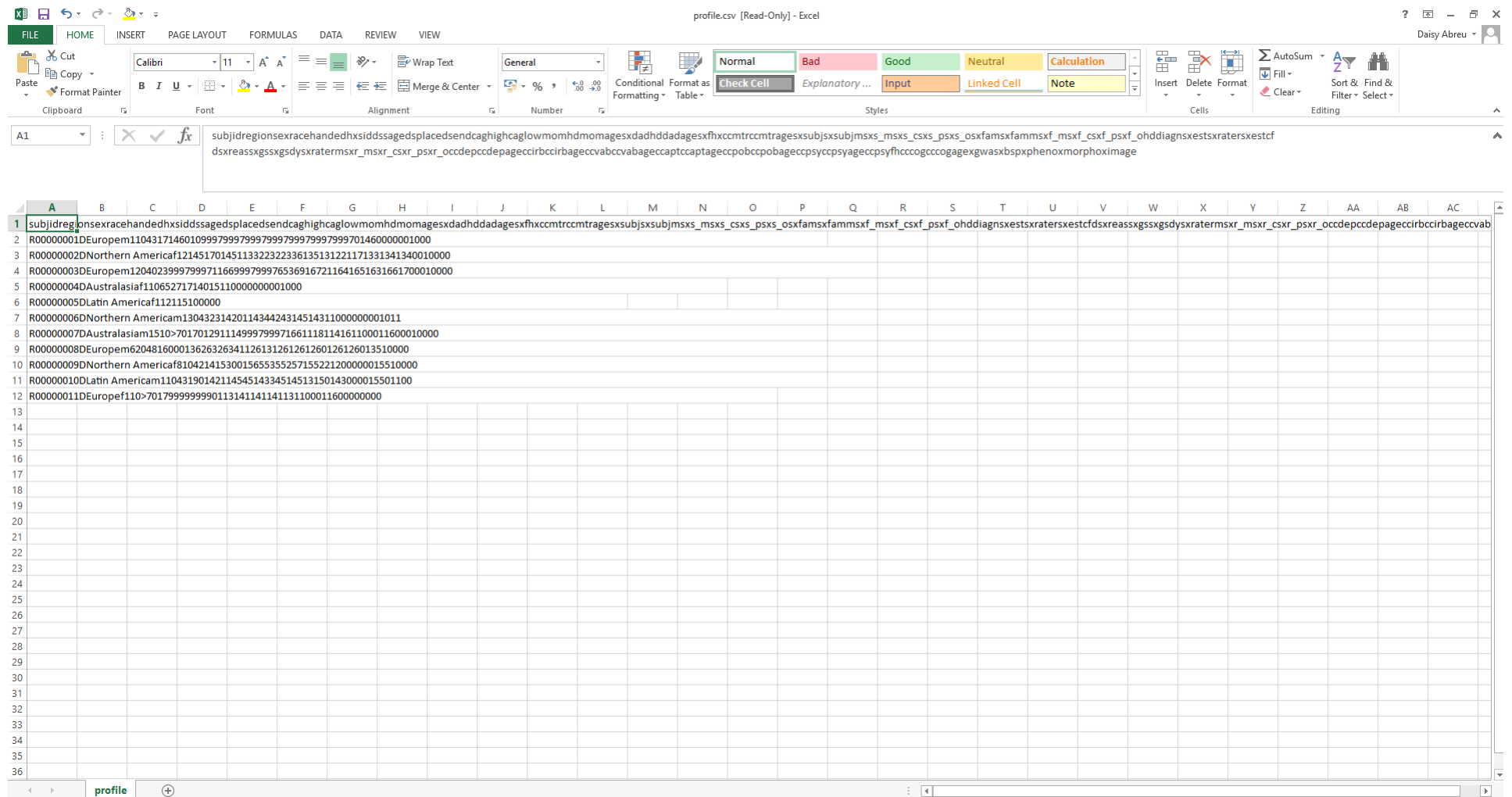
```
load("Rdata_directory")
```

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





profile.csv [Read-Only] - Excel

FILE HOME INSERT PAGE LAYOUT FORMULAS DATA REVIEW VIEW

From Access From Web From Text From Other Sources Get External Data Existing Connections Refresh All Edit Links Connections Sort & Filter Filter Clear Reapply Advanced Text to Columns Flash Fill Remove Duplicates Data Validation Consolidate What-If Analysis Relationships Group Ungroup Subtotal Outline Show Detail Hide Detail

A1 : subjldregionsexracehandedhxsiddssagedsplacedsendcaghighcaglowmomhdmomagesxdadhdadagesfxhccmtrcmtragesxsubjsxsubjmsxs\_msxs\_csxs\_pxsx\_oxsfamsxfammsxf\_msxf\_csxf\_psf\_ohddiagnskestxstratersxestcfdsreassxgssxgdsyratermsxr\_msxr\_csxr\_psr\_ocdepccdepagccirbcbirbageccvabccvabageccaptcaptageccpccpobageccpsycpsysageccpsyfhhccogccogagexgwaxbspxphenoxmorphoximage

Subjld	Region	Sex	Race	Handed	Hxsid	Issage	Spplace
R0000001D	Europe	M	W	R	0		
R0000002D	Northern America	M	W	R	1		
R0000003D	Europe	M	W	R	0		
R0000004D	Australasia	M	W	R	0	65	2

Convert Text to Columns Wizard - Step 2 of 3

This screen lets you set the delimiters your data contains. You can see how your text is affected in the preview below.

Delimiters

Tab

Semicolon  Treat consecutive delimiters as one

Comma

Space

Other:

Text qualifier:

Data preview

Subjld	Region	Sex	Race	Handed	Hxsid	Issage	Spplace
R0000001D	Europe	M	W	R	0		
R0000002D	Northern America	M	W	R	1		
R0000003D	Europe	M	W	R	0		
R0000004D	Australasia	M	W	R	0	65	2

Cancel < Back Next > Finish

**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”.

The screenshot shows the 'Convert Text to Columns Wizard - Step 3 of 3' dialog box in Microsoft Excel. The dialog is titled 'Convert Text to Columns Wizard - Step 3 of 3'. It contains the following elements:

- Column data format:** Radio buttons for 'General', 'Text', 'Date: DMY', and 'Do not import column (skip)'. The 'Text' option is selected.
- Destination:** A text box containing 'SA\$1'.
- Data preview:** A table showing the first few rows of data. The columns are: 'Region', 'Sex', 'Race', 'Handed', 'Hxsid', 'Hsage', and 'Hsplace'. The data rows are:
 

Region	Sex	Race	Handed	Hxsid	Hsage	Hsplace
Europe	M	L	L	0		
Northern America	F	L	2	0		
Europe	F	L	2	0		
Australasia	F	L	1	0	65	2

Buttons at the bottom include 'Cancel', '< Back', 'Next >', and 'Finish'.

NB: The data files *pharmacotx* and *nutsuppl* contain two columns ‘cmtrt\_decod’ and ‘cmdstot’ 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.



# Enroll-HD: Download and Import Data | 2022-05-25



profile.csv - Excel

Microsoft Excel ribbon: FILE, HOME, INSERT, PAGE LAYOUT, FORMULAS, DATA, REVIEW, VIEW. Data Tools group includes: Connections, Sort & Filter, Text to Columns, Flash Fill, Remove Duplicates, Data Validation, Consolidate, What-If Analysis, Relationships, Group, Ungroup, Subtotal, Outline.

Formula bar: A1 : subjid

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	
1	subjid	region	sex	race	handed	hxsid	dssage	dsplace	dsend	caghigh	caglow	momhd	momages	dadhd	dadagesx	fhx	ccmtr	ccmtrage	sxsubj	sxsubjm	sxs_m	sxs_c	sxs_p	sxs_o	sxfam	sxfamm	sxf_m	sxf_c	sxf_p	sx
2	R0000000	Europe	m		1	1	0			43	17	1	46	0		1	0		9997	9997					9997	9997				
3	R0000000	Northern	f		1	2	1			45	17	0		1	45	1	1	33	22	3					22	3				
4	R0000000	Europe	m		1	2	0			40	23	9997		9997		1	1	66	9997	9997					65	3				
5	R0000000	Australasi	f		1	1	0	65	2	7	17	14	0		1	51	1	0												
6	R0000000	Latin Ame	f		1	1				21	15					1														
7	R0000000	Northern	m		1	3	0			43	23	1	42	0		1	1	43	44	2					43	1				
8	R0000000	Australasi	m		15	1	0			>70	17	0		1	29	1	1	14	9997	9997					16	6	1	1		
9	R0000000	Europe	m		6	2	0			48	16	0		0		0	1	36	26	3					26	3				
10	R0000000	Northern	f		8	1	0			42	14	1	53	0		0	1	56	55	3					55	2				
11	R0000001	Latin Ame	m		1	1	0			43	19	0		1	42	1	1	45	45	1					43	3				
12	R0000001	Europe	f		1	1	0			>70	17	9999		9999		0	1	13	14	1					14	1				