

Reading & inspecting data in R

- R can read and write data from a multitude of different sources
 - Text (csv, tsv, ...)
 - Stata
 - SPSS
 - SAS
 - Excel
 - ...
- And using lots of different packages and functions

Data type	Extension	Function	Package
Comma separated values	csv	<code>read.csv()</code>	utils (default)
		<code>read_csv()</code>	readr (tidyverse)
Tab separated values	tsv	<code>read_tsv()</code>	readr
Other delimited formats	txt	<code>read.table()</code>	utils
		<code>read_table()</code>	readr
		<code>read_delim()</code>	readr
Stata version 13-14	dta	<code>readdta()</code>	haven
Stata version 7-12	dta	<code>read.dta()</code>	foreign
SPSS	sav	<code>read.spss()</code>	foreign
SAS	sas7bdat	<code>read.sas7bdat()</code>	sas7bdat
Excel	xlsx, xls	<code>read_excel()</code>	readxl (tidyverse)

raw -> technically correct data

- `read.table` and its cousins (`utils` package)
 - `read.csv`: *comma* separated values with *period* as decimal separator.
 - `read.csv2`: *semicolon* separated values with *comma* as decimal separator.
 - `read.delim`: *tab-delimited* files with *period* as decimal separator.
 - `read.delim2`: *tab-delimited* files with *comma* as decimal separator.
 - `read.fwf`: data with a predetermined number of bytes per column.

Example

- Simple base R option

```
surveys <- read.csv("data/portal_data_joined.csv")
```

Downloading from the web

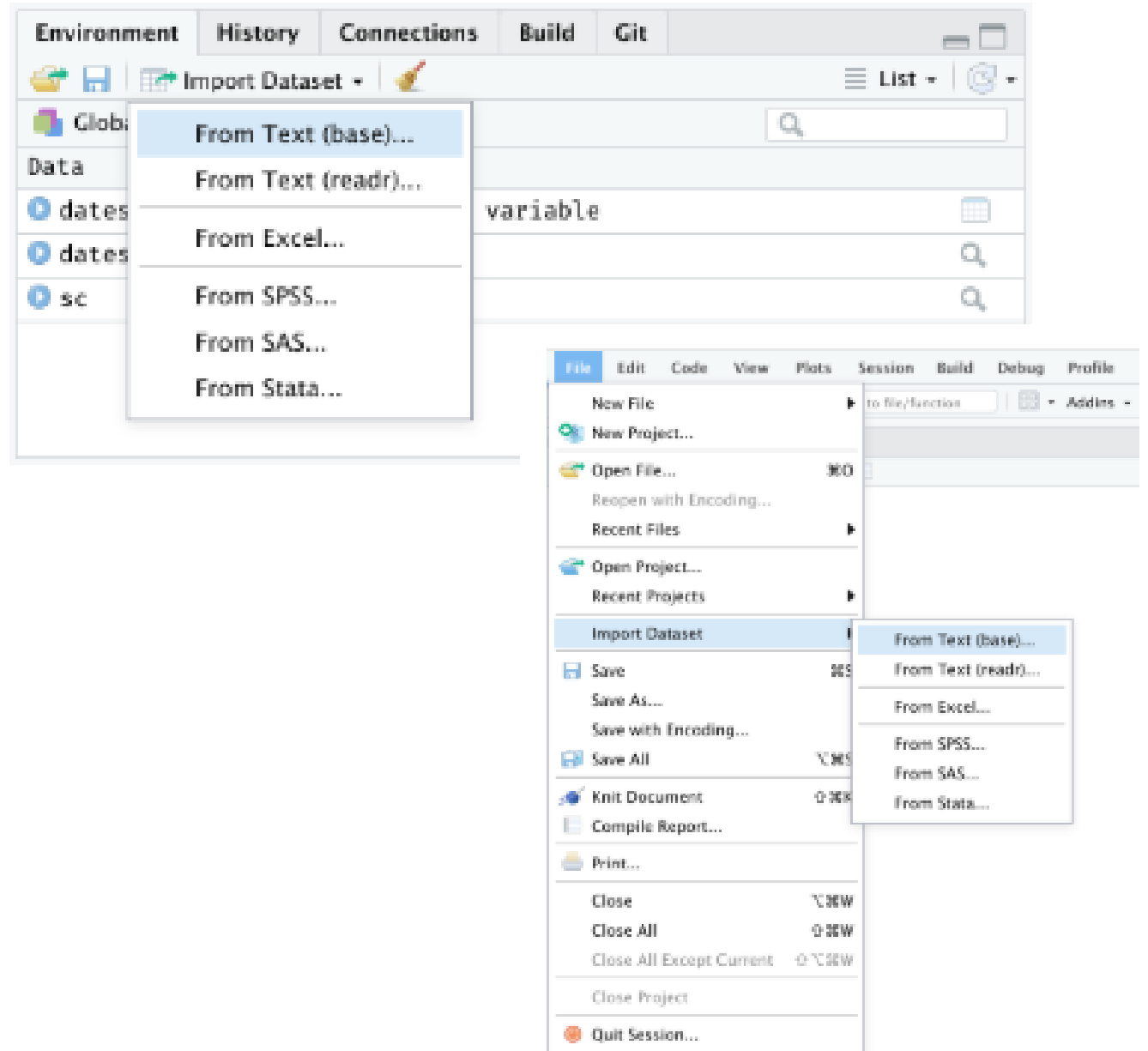
- We can download data before reading it in to R
- A common function for this is `download.file()`
- E.g.

```
download.file(url="https://ndownloader.figshare.com/files/22921  
             destfile = "data/portal_data_joined.csv")
```

- Can download from e.g. GitHub, figshare, googlesdrive, ...

Via GUI

- 3 categories
 - Text data
 - Excel data
 - Statistical data
- Options
 - Environment pane
 - File menu



Importing using readr

- Provides support for
 - Import from the file system or a url
 - Change column data types
 - Skip or include-only columns
 - Rename the data set
 - Skip the first N rows
 - Use the header row for column names
 - Trim spaces in names
 - Change the column delimiter
 - Encoding selection
 - Select quote, escape, comment and NA identifiers

Using stringsAsFactors=FALSE

- By default, when building or importing a data frame, the columns that contain characters (i.e. text) are coerced (= converted) into factors. Depending on what you want to do with the data, you may want to keep these columns as character. To do so, `read.csv()` and `read.table()` have an argument called `stringsAsFactors` which can be set to `FALSE`.
- In most cases, it is preferable to set `stringsAsFactors = FALSE` when importing data and to convert as a factor only the columns that require this data type.

Importing data from Text files

Import Dataset

Name:

Encoding:

Heading: ☒ Yes ☐ No

Row names:

Separator:

Decimal:

Quote:

Comment:

na.strings:

☒ Strings as factors

Input File

```
Sepal.Length,Sepal.Width,Petal.Length,Petal.Width,Species
5.1,3.5,1.4,0.2,setosa
4.9,3.1,1.4,0.2,setosa
4.7,3.2,1.3,0.2,setosa
4.6,3.1,1.5,0.2,setosa
5.3,6.1,4.0,2.0,setosa
5.4,3.9,1.7,0.4,setosa
4.6,3.4,1.4,0.3,setosa
5.3,4.1,5.0,2.0,setosa
4.4,2.9,1.4,0.2,setosa
4.9,3.1,1.5,0.1,setosa
5.4,3.7,1.5,0.2,setosa
4.8,3.4,1.6,0.2,setosa
4.8,3.1,1.4,0.1,setosa
```

Data Frame

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
5.4	3.9	1.7	0.4	setosa
4.6	3.4	1.4	0.3	setosa
5.0	3.4	1.5	0.2	setosa
4.4	2.9	1.4	0.2	setosa
4.9	3.1	1.5	0.1	setosa
5.4	3.7	1.5	0.2	setosa
4.8	3.4	1.6	0.2	setosa
4.8	3.0	1.4	0.1	setosa

Import Cancel

Importing data from Excel files

- The Excel importer provides support to:
 - Import from the file system or a url
 - Change column data types
 - Skip columns
 - Rename the data set
 - Select an specific Excel sheet
 - Skip the first N rows
 - Select NA identifiers

Import Excel Data

File/URL:

<http://www.fns.usda.gov/sites/default/files/pd/sisummar.xls>

Update

Data Preview:

NATIONAL SCHOOL LUNCH PROGRAM PARTICIPATION AND LUNCHES SERVED	X_1	X_2	X_3	X_4	X_5	X_6
(character)	(character)	(character)	(character)	(character)	(character)	(character)
(Data as of October 05, 2018)	NA	NA	NA	NA	NA	NA
	-----Average Participation-----	NA	NA	NA	NA	NA
Fiscal		Reduced	Full		Total Lunches	Percent Free/RP
Year	Free	Price	Price	Total	Served	of Total
	-----Millions-----	NA	NA	NA	NA	%
1969	2.8999999999999999	1}	38.1	19.199999999999999	3368.1999999999998	15.1
1970	4.5999999999999996	1}	17.600000000000001	22.399999999999999	3565.0999999999999	20.699999999999999
1971	5.7999999999999998	0.5	17.600000000000001	24.399999999999999	3848.3999999999998	26.399999999999999
1972	7.2999999999999998	0.5	16.600000000000001	24.399999999999999	3972.0999999999999	32.399999999999999
...

Previewing first 50 entries.

Import Options:

Name: Max Rows: ☒ First Row as Names
 Sheet: Skip: ☒ Open Data Viewer
 Range: NA:

Code Preview:

```
library(readxl)
url <- "http://www.fns.usda.gov/sites/default/files/pd/sisummar.xls"
destfile <- "sisummar.xls"
curl::curl_download(url, destfile)
sisummar <- read_excel(destfile)
View(sisummar)
```

Reading Excel files using readxl

Import

Cancel

- We can clean this up by
 - skipping 6 rows from this file
 - unchecking the "First Row as Names" checkbox.
- The file is looking better but some columns are being displayed as strings when they are clearly numerical data.
- We can fix this by selecting "numeric" from the column dropdown.
- The final step is to click "Import" to run the code under "Code Preview" and import the data into RStudio, the final result should look as follows:

.RData, .Rda, .Rds

- Store R objects with properties
- .Rda is just a short name for .Rdata
 - `save()`, `load()`, `attach()`,...
- .Rds use to restore R object with different name
 - `readRDS()`, `saveRDS()`,...

Inspecting data structure

- Once data is read-in there are several way to view and interrogate them
- The most simple is to type the name of the variable
- Can you see any problems with this?

Other options

- View the top of the data
- `head(<name>)`
- Equivalently we can use
- `tail(<name>)`

Compact display of contents

- The head and tail are limited in terms of information
- A good overall summary in base R is
- `str(<name>)`

```
str(metadata)

'data.frame':  12 obs. of  3 variables:
 $ genotype : Factor w/ 2 levels "KO","Wt": 2 2 2 1 1 1 2 2 2 1 ...
 $ celltype  : Factor w/ 2 levels "typeA","typeB": 1 1 1 1 1 1 2 2 2 2 ...
 $ replicate: num  1 2 3 1 2 3 1 2 3 1 ...
```

List of some data inspection functions

- All data structures - content display:
 - `str()` : compact display of data contents (env.)
 - `class()` : data type (e.g. character, numeric, etc.) of vectors and data structure of dataframes, matrices, and lists.
 - `summary()` : detailed display, including descriptive statistics, frequencies
 - `head()` : will print the beginning entries for the variable
 - `tail()` : will print the end entries for the variable
- Vector and factor variables:
 - `length()` : returns the number of elements in the vector or factor
- Dataframe and matrix variables:
 - `dim()` : returns dimensions of the dataset
 - `nrow()` : returns the number of rows in the dataset
 - `ncol()` : returns the number of columns in the dataset
 - `rownames()` : returns the row names in the dataset
 - `colnames()` : returns the column names in the dataset

Viewing data

- You can view the data (when data frame, matrix etc) in the data panel in several ways:

`View (<name>)`

`Edit (<name>)`

- Left click with the mouse on the name of the variable in the Environment tab
- Hover over the variable name in the console or scripts and press F2

Writing to Excel

- Format, styling and editing inside of R
- `library(openxlsx)`
 - `l <- list(iris = iris, mtcars = mtcars, chickwts = chickwts, quakes = quakes)`
 - `openxlsx::write.xlsx(l, file="datasets.xlsx")`
- Much more fancy features
- E.g.
- <https://cran.r-project.org/web/packages/openxlsx/vignettes/formatting.pdf>

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do									
<div> <div> <div>Cut Copy Format Painter</div> <div>Clipboard</div> </div> <div> <div>Calibri 11 A A</div> <div>B I U</div> <div>Font</div> </div> <div> <div> <div> <div></div> <div></div> </div> <div> <div></div> <div></div> </div> </div> <div> <div>Wrap Text</div> <div>Merge & Center</div> <div>Alignment</div> </div> </div> </div>									
1									
A	B	C	D	E	F	G	H	I	J
	Date	Logical	Curren	Accoun	hLink	Percen	TinyNu		
	#####	TRUE	-\$2.00	-£ 2.00	https://CF	-100.00%	5.66E-10		
	#####	FALSE	-\$1.00	-£ 1.00	https://CF	-50.00%	8.55E-10		
	#####	TRUE	\$0.00	£ -	https://CF	0.00%	5.80E-10		
	#####	TRUE	\$1.00	£ 1.00	https://CF	50.00%	1.52E-10		
	#####	FALSE	\$2.00	£ 2.00	https://CF	100.00%	3.75E-10		
	Date	Logical	Curren	Accoun	hLink	Percen	TinyNu		
	#####	TRUE	-\$2.00	-£ 2.00	https://CF	-100.00%	5.66E-10		
	#####	FALSE	-\$1.00	-£ 1.00	https://CF	-50.00%	8.55E-10		
	#####	TRUE	\$0.00	£ -	https://CF	0.00%	5.80E-10		
	#####	TRUE	\$1.00	£ 1.00	https://CF	50.00%	1.52E-10		
	#####	FALSE	\$2.00	£ 2.00	https://CF	100.00%	3.75E-10		
	Date	Logical	Curren	Accoun	hLink	Percen	TinyNu	per	
	#####	TRUE	-\$2.00	-£ 2.00	https://CF	-100.00%	5.66E-10		
	#####	FALSE	-\$1.00	-£ 1.00	https://CF	-50.00%	8.55E-10		
	#####	TRUE	\$0.00	£ -	https://CF	0.00%	5.80E-10		
	#####	TRUE	\$1.00	£ 1.00	https://CF	50.00%	1.52E-10		
	#####	FALSE	\$2.00	£ 2.00	https://CF	100.00%	3.75E-10		
	Date	Logical	Curren	Accoun	hLink	Percen	TinyNu		
	#####	TRUE	-\$2.00	-£ 2.00	https://CF	-100.00%	5.66E-10		
	#####	FALSE	-\$1.00	-£ 1.00	https://CF	-50.00%	8.55E-10		
	#####	TRUE	\$0.00	£ -	https://CF	0.00%	5.80E-10		
	#####	TRUE	\$1.00	£ 1.00	https://CF	50.00%	1.52E-10		
	#####	FALSE	\$2.00	£ 2.00	https://CF	100.00%	3.75E-10		