

# Statistical Moments and Beyond

# Outline: Statistical Moments and Beyond

- What is statistics, and what are the importance of statistics?
- Data exploration, manipulation.
- Examining distributions
- Basic descriptive data visualisation



# An introduction to *tidy data management* **with R**

# {Tidy} data management

## An Introduction and The Why

# {Tidy} data management: The Why



# {Tidy} data management: The Why

O1_LONGITUDE	O2_LATITUDE	One of main goals in life has been to make my parents proud	Feeling of happiness	State of health (subjective)	How often do you pray	Year of birth	Age recoded (6 intervals)	Employment status	country code
-1.537885	55.078311	Agree strongly	Very happy	Very good	Never, practically never	2000	16-24	Part time (less than 30 hours a week)	GBR
-1.532713	55.074637	Agree	Quite happy	Poor	Once a day	1950	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.608998	55.132695	Agree	Very happy	Poor	Never, practically never	1971	45-54	Part time (less than 30 hours a week)	GBR
-1.608338	55.1344	Don't know	Very happy	Poor	Several times each week	1988	25-34	Unemployed	GBR
-1.602372	55.128,665	Agree	Quite happy	Fair	Never, practically never	1966	55-64	Self employed	GBR
-1.603687	55.12734	Agree	Very happy	Good	Never, practically never	1947	65 and more years	Retired/pensioned	GBR
-1.601168	55.134646	Disagree	Not very happy	Good	Once a day	1944	65 and more years	Retired/pensioned	GBR
-1.602418	55.130346	Agree	Quite happy	Fair	Never, practically never	1954	65 and more years	Retired/pensioned	GBR
-1.608852	55.13645	Agree	Very happy	Good	Several times each week	1949	65 and more years	Retired/pensioned	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.458406	55.046.29	Disagree	Quite happy		Several times each week	1937	65 and more years	Retired/pensioned	GBR
-1.461515	55.051.383	Agree	Not very happy		Several times each week	1960	55-64	Retired/pensioned	GBR
-1.465077	55.051893	Agree	N/A	Good	Less often	1947	N/A	Retired/pensioned	GBR
-1.467017	55.054274	Agree	Quite happy	Good	Never, practically never	1978	35-44	Part time (less than 30 hours a week)	GBR
-1.464305	55.054585	Agree strongly	Very happy	Very good	Never, practically never	1951	65 and more years	Retired/pensioned	GBR
-1.461513	55.044554	Strongly disagree	Quite happy	Good	Less often	1982	35-44	Full time (30 hours a week or more)	GBR
-1.452616	55.043975	NA	Very happy	Good	Less often	1939	65+	Retired/pensioned	GBR
-1.454677	55.043175	Agree	Very happy	Good	Never, practically never	1950	65 and more years	Retired/pensioned	GBR
-1.505645	54.995998	Agree	Very happy	Fair	Never, practically never	1906	65 and more years	Retired/pensioned	GBR

An irrelevant column

# {Tidy} data management: The Why

O1_LONGITUDE	O2_LATITUDE	One of main goals in life has been to make my parents proud	Feeling of happiness	State of health (subjective)	How often do you pray	Year of birth	Age recoded (6 intervals)	Employment status	country code
-1.537885	55.078311	Agree strongly	Very happy	Very good	Never, practically never	2000	16-24	Part time (less than 30 hours a week)	GBR
-1.532713	55.074637	Agree	Quite happy	Poor	Once a day	1950	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.608998	55.132695	Agree	Very happy	Poor	Never, practically never	1971	45-54	Part time (less than 30 hours a week)	GBR
-1.608338	55.1344	Don't know	Very happy	Poor	Several times each week	1988	25-34	Unemployed	GBR
-1.602372	55.128.665	Agree	Quite happy	Fair	Never, practically never	1966	55-64	Self employed	GBR
-1.603687	55.12734	Agree	Very happy	Good	Never, practically never	1947	65 and more years	Retired/pensioned	GBR
-1.601168	55.134646	Disagree	Not very happy	Good	Once a day	1944	65 and more years	Retired/pensioned	GBR
-1.602418	55.130346	Agree	Quite happy	Fair	Never, practically never	1954	65 and more years	Retired/pensioned	GBR
-1.608852	55.13645	Agree	Very happy	Good	Several times each week	1949	65 and more years	Retired/pensioned	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.458406	55.046.29	Disagree	Quite happy		Several times each week	1937	65 and more years	Retired/pensioned	GBR
-1.461515	55.051.383	Agree	Not very happy		Several times each week	1960	55-64	Retired/pensioned	GBR
-1.465077	55.051893	Agree	N/A	Good	Less often	1947	N/A	Retired/pensioned	GBR
-1.467017	55.054274	Agree	Quite happy	Good	Never, practically never	1978	35-44	Part time (less than 30 hours a week)	GBR
-1.464305	55.054585	Agree strongly	Very happy	Very good	Never, practically never	1951	65 and more years	Retired/pensioned	GBR
-1.461513	55.044554	Strongly disagree	Quite happy	Good	Less often	1982	35-44	Full time (30 hours a week or more)	GBR
-1.452616	55.043975	NA	Very happy	Good	Less often	1939	65+	Retired/pensioned	GBR
-1.454677	55.043175	Agree	Very happy	Good	Never, practically never	1950	65 and more years	Retired/pensioned	GBR
-1.505645	54.995998	Agree	Very happy	Fair	Never, practically never	1906	65 and more years	Retired/pensioned	GBR

Duplicate observations

# {Tidy} data management: The Why

O1_LONGITUDE	O2_LATITUDE	One of main goals in life has been to make my parents proud	Feeling of happiness	State of health (subjective)	How often do you pray	Year of birth	Age recoded (6 intervals)	Employment status	country code
-1.537885	55.078311	Agree strongly	Very happy	Very good	Never, practically never	2000	16-24	Part time (less than 30 hours a week)	GBR
-1.532713	55.074637	Agree	Quite happy	Poor	Once a day	1950	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.608998	55.132695	Agree	Very happy	Poor	Never, practically never	1971	45-54	Part time (less than 30 hours a week)	GBR
-1.608338	55.1344	Don't know	Very happy	Poor	Several times each week	1988	25-34	Unemployed	GBR
-1.602372	55.128,665	Agree	Quite happy	Fair	Never, practically never	1966	55-64	Self employed	GBR
-1.603687	55.12734	Agree	Very happy	Good	Never, practically never	1947	65 and more years	Retired/pensioned	GBR
-1.601168	55.134646	Disagree	Not very happy	Good	Once a day	1944	65 and more years	Retired/pensioned	GBR
-1.602418	55.130346	Agree	Quite happy	Fair	Never, practically never	1954	65 and more years	Retired/pensioned	GBR
-1.608852	55.13645	Agree	Very happy	Good	Several times each week	1949	65 and more years	Retired/pensioned	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.458406	55.046.29	Disagree	Quite happy		Several times each week	1937	65 and more years	Retired/pensioned	GBR
-1.461515	55.051.383	Agree	Not very happy		Several times each week	1960	55-64	Retired/pensioned	GBR
-1.465077	55.051893	Agree	N/A	Good	Less often	1947	N/A	Retired/pensioned	GBR
-1.467017	55.054274	Agree	Quite happy	Good	Never, practically never	1978	35-44	Part time (less than 30 hours a week)	GBR
-1.464305	55.054585	Agree strongly	Very happy	Very good	Never, practically never	1951	65 and more years	Retired/pensioned	GBR
-1.461513	55.044554	Strongly disagree	Quite happy	Good	Less often	1982	35-44	Full time (30 hours a week or more)	GBR
-1.452616	55.043975	NA	Very happy	Good	Less often	1939	65+	Retired/pensioned	GBR
-1.454677	55.043175	Agree	Very happy	Good	Never, practically never	1950	65 and more years	Retired/pensioned	GBR
-1.505645	54.995998	Agree	Very happy	Fair	Never, practically never	1906	65 and more years	Retired/pensioned	GBR

Outliers

# {Tidy} data management: The Why

Structural  
Errors

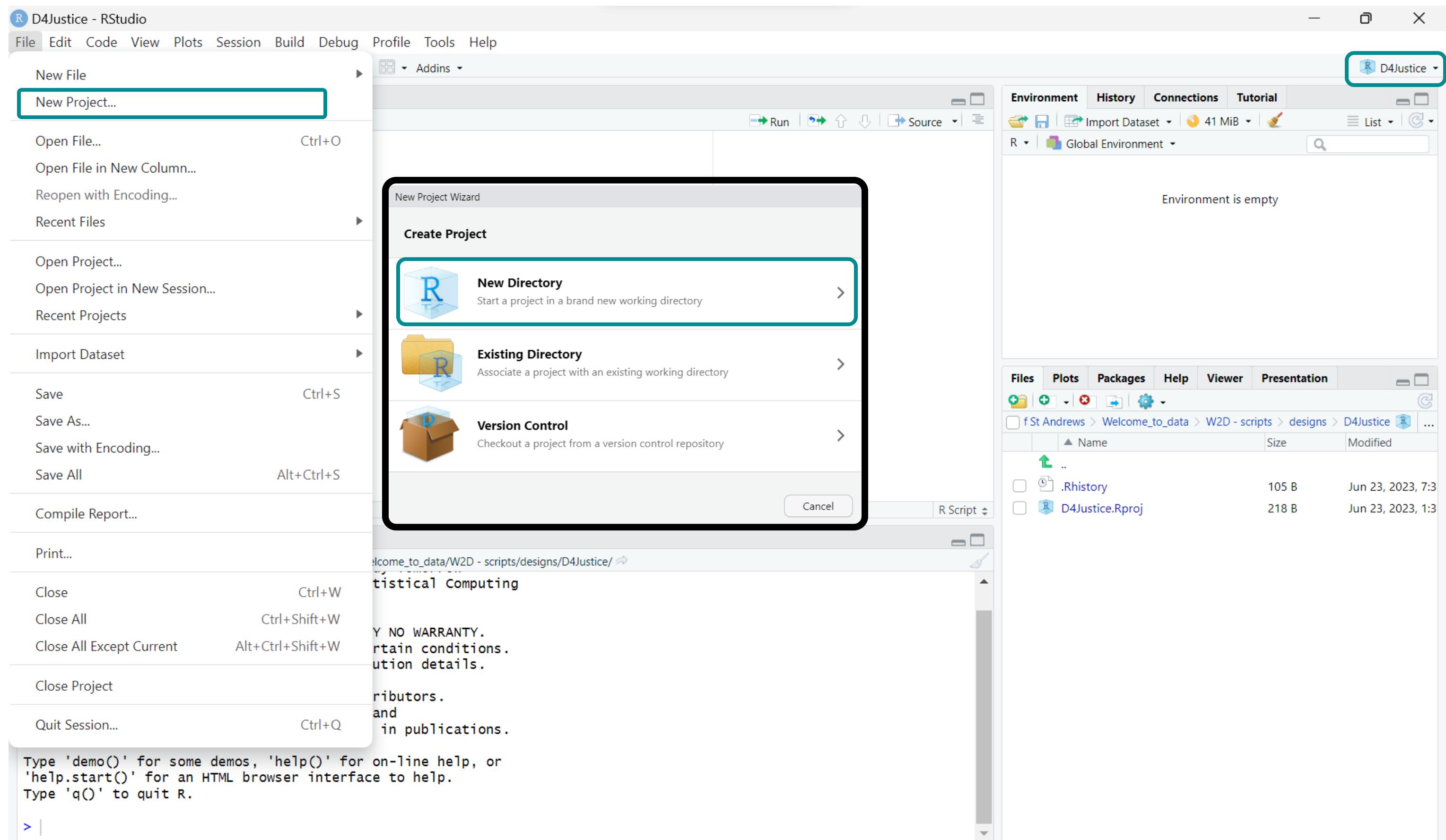
O1_LONGITUDE	O2_LATITUDE	One of main goals in life has been to make my parents proud	Feeling of happiness	State of health (subjective)	How often do you pray	Year of birth	Age recoded (6 intervals)	Employment status	country code
-1.537885	55.078311	Agree strongly	Very happy	Very good	Never, practically never	2000	16-24	Part time (less than 30 hours a week)	GBR
-1.532713	55.074637	Agree	Quite happy	Poor	Once a day	1950	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.608998	55.132695	Agree	Very happy	Poor	Never, practically never	1971	45-54	Part time (less than 30 hours a week)	GBR
-1.608338	55.1344	Don't know	Very happy	Poor	Several times each week	1988	25-34	Unemployed	GBR
-1.602372	55.128,665	Agree	Quite happy	Fair	Never, practically never	1966	55-64	Self employed	GBR
-1.603687	55.12734	Agree	Very happy	Good	Never, practically never	1947	65 and more years	Retired/pensioned	GBR
-1.601168	55.134646	Disagree	Not very happy	Good	Once a day	1944	65 and more years	Retired/pensioned	GBR
-1.602418	55.130346	Agree	Quite happy	Fair	Never, practically never	1954	65 and more years	Retired/pensioned	GBR
-1.608852	55.13645	Agree	Very happy	Good	Several times each week	1949	65 and more years	Retired/pensioned	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.458406	55,046.29	Disagree	Quite happy		Several times each week	1937	65 and more years	Retired/pensioned	GBR
-1.461515	55.051.383	Agree	Not very happy		Several times each week	1960	55-64	Retired/pensioned	GBR
-1.465077	55.051893	Agree	N/A	Good	Less often	1947	N/A	Retired/pensioned	GBR
-1.467017	55.054274	Agree	Quite happy	Good	Never, practically never	1978	35-44	Part time (less than 30 hours a week)	GBR
-1.464305	55.054585	Agree strongly	Very happy	Very good	Never, practically never	1951	65 and more years	Retired/pensioned	GBR
-1.461513	55.044554	Strongly disagree	Quite happy	Good	Less often	1982	35-44	Full time (30 hours a week or more)	GBR
-1.452616	55.043975	NA	Very happy	Good	Less often	1939	65+	Retired/pensioned	GBR
-1.454677	55.043175	Agree	Very happy	Good	Never, practically never	1950	65 and more years	Retired/pensioned	GBR
-1.505645	54.995998	Agree	Very happy	Fair	Never, practically never	1906	65 and more years	Retired/pensioned	GBR

# {Tidy} data management: The Why

Missing  
Values

O1_LONGITUDE	O2_LATITUDE	One of main goals in life has been to make my parents proud	Feeling of happiness	State of health (subjective)	How often do you pray	Year of birth	Age recoded (6 intervals)	Employment status	country code
-1.537885	55.078311	Agree strongly	Very happy	Very good	Never, practically never	2000	16-24	Part time (less than 30 hours a week)	GBR
-1.532713	55.074637	Agree	Quite happy	Poor	Once a day	1950	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.535722	55.073129	Agree	Quite happy	Good	Several times each week	1952	65 and more years	Retired/pensioned	GBR
-1.608998	55.132695	Agree	Very happy	Poor	Never, practically never	1971	45-54	Part time (less than 30 hours a week)	GBR
-1.608338	55.1344	Don't know	Very happy	Poor	Several times each week	1988	25-34	Unemployed	GBR
-1.602372	55.128,665	Agree	Quite happy	Fair	Never, practically never	1966	55-64	Self employed	GBR
-1.603687	55.12734	Agree	Very happy	Good	Never, practically never	1947	65 and more years	Retired/pensioned	GBR
-1.601168	55.134646	Disagree	Not very happy	Good	Once a day	1944	65 and more years	Retired/pensioned	GBR
-1.602418	55.130346	Agree	Quite happy	Fair	Never, practically never	1954	65 and more years	Retired/pensioned	GBR
-1.608852	55.13645	Agree	Very happy	Good	Several times each week	1949	65 and more years	Retired/pensioned	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.200647	54.688283	Agree	Quite happy	Poor	Never, practically never	1993	25-34	Full time (30 hours a week or more)	GBR
-1.458406	55.046.29	Disagree	Quite happy		Several times each week	1937	65 and more years	Retired/pensioned	GBR
-1.461515	55.051.383	Agree	Not very happy		Several times each week	1960	55-64	Retired/pensioned	GBR
-1.465077	55.051893	Agree	N/A	Good	Less often	1947	N/A	Retired/pensioned	GBR
-1.467017	55.054274	Agree	Quite happy	Good	Never, practically never	1978	35-44	Part time (less than 30 hours a week)	GBR
-1.464305	55.054585	Agree strongly	Very happy	Very good	Never, practically never	1951	65 and more years	Retired/pensioned	GBR
-1.461513	55.044554	Strongly disagree	Quite happy	Good	Less often	1982	35-44	Full time (30 hours a week or more)	GBR
-1.452616	55.043975	NA	Very happy	Good	Less often	1939	65+	Retired/pensioned	GBR
-1.454677	55.043175	Agree	Very happy	Good	Never, practically never	1950	65 and more years	Retired/pensioned	GBR
-1.505645	54.995998	Agree	Very happy	Fair	Never, practically never	1906	65 and more years	Retired/pensioned	GBR

# {Tidy} data management: The Why



# {Tidy} data management

## The How

# Tidy data management: The How

The screenshot shows the RStudio interface with the following components:

- Code Editor:** Displays an R script named "Untitled1" containing the following code:

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
```
- Environment View:** Shows the Global Environment tab with the message "Environment is empty".
- File Explorer:** Shows the project structure under "University of St Andrews > Welcome\_to\_data > W2D - scripts > designs > D4Justice". It lists two files: ".Rhistory" (105 B, Jun 24, 2023, 9:24 PM) and "D4Justice.Rproj" (218 B, Nov 6, 2023, 4:24 PM).
- Console:** Displays the R startup message and the command prompt "> |".

# Tidy data management: The How

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
7
8 participants <- data.frame (age, gender,
9                               education, country)
```

# Tidy data management: The How

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
7
8 participants <- data.frame (age, gender,
9                               education, country)
10
11 dim (participants)
```

[1] 5 4

# Tidy data management: The How

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
7
8 participants <- data.frame (age, gender,
9                               education, country)
10
11 head (participants, 3)
```

	age	gender	education	country
1	34	M	Tertiary	Nigeria
2	23	M	Primary	Britain
3	56	F	Primary	Peru

# Tidy data management: The How

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
7
8 participants <- data.frame (age, gender,
9                               education, country)
10
11 head (participants, 3)
```

```
age gender education country
1 34      M   Tertiary Nigeria
2 23      M   Primary  Britain
3 56      F   Primary    Peru
```

```
1 tail (participants, 3)
```

```
age gender education country
3 56      F   Primary    Peru
4 34      F Secondary   USA
5 46      M   Tertiary  <NA>
```

# Tidy data management: The How

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
7
8 participants <- data.frame (age, gender,
9                               education, country)
10
11 head (participants, 3)
```

```
age gender education country
1 34      M   Tertiary Nigeria
2 23      M   Primary  Britain
3 56      F   Primary    Peru
```

```
1 tail (participants, 3)
```

```
age gender education country
3 56      F   Primary    Peru
4 34      F   Secondary   USA
5 46      M   Tertiary  <NA>
```

```
1 str (participants)
```

```
'data.frame': 5 obs. of 4 variables:
 $ age     : num 34 23 56 34 46
 $ gender   : chr "M" "M" "F" "F" ...
 $ education: chr "Tertiary" "Primary" "Primary" "Secondary"
 ...
 $ country  : chr "Nigeria" "Britain" "Peru" "USA" ...
```

# Tidy data management: The How

```
1 age <- c(34, 23, 56, 34, 46)
2 gender <- c("M", "M", "F", "F", "M")
3 education <- c("Tertiary", "Primary",
4                 "Primary", "Secondary", "Tertiary")
5 country <- c("Nigeria", "Britain", "Peru",
6                  "USA", NA)
7
8 participants <- data.frame (age, gender,
9                               education, country)
10
11 head (participants, 3)
```

```
age gender education country
1 34      M   Tertiary Nigeria
2 23      M   Primary  Britain
3 56      F   Primary    Peru
```

```
1 tail (participants, 3)
```

```
age gender education country
3 56      F   Primary    Peru
4 34      F   Secondary   USA
5 46      M   Tertiary  <NA>
```

```
1 str (participants)
```

```
'data.frame': 5 obs. of 4 variables:
 $ age     : num 34 23 56 34 46
 $ gender   : chr "M" "M" "F" "F" ...
 $ education: chr "Tertiary" "Primary" "Primary" "Secondary"
 ...
 $ country  : chr "Nigeria" "Britain" "Peru" "USA" ...
```

```
1 str (participants$age)
```

```
num [1:5] 34 23 56 34 46
```

# {Tidy} data management

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

# {Tidy} data management:

## A case study using the **World Value Survey**

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

The screenshot shows the WVS website's "Data Download" section for Wave 7 (2017-2022). The left sidebar has a tree view of the survey's structure. The main content area is organized into several sections:

- United Kingdom - Great Britain 2022**
- Questionnaire**: WVS7 Questionnaire Great Britain 2022 English.pdf
- Sampling & Methodology**: WVS7 Methodology Report Great Britain.pdf, WVS7 Sample Design Great Britain 2022.pdf, WVS7 Information about the team Great Britain 2022.pdf
- Codebook & Results**: World Values Survey Wave 7 (2017-2022) UK - Great Britain v5.0
- Data Files**: WVS Wave 7 UK - Great Britain Csv v5.0, WVS Wave 7 UK - Great Britain CsvTxt v5.0, WVS Wave 7 UK - Great Britain Excel v5.0, WVS Wave 7 UK - Great Britain ExcelTxt v5.0, WVS Wave 7 UK - Great Britain Spss v5.0, WVS Wave 7 UK - Great Britain Stata v5.0

# {Tidy} data management:

## A case study using the **World Value Survey**

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate



# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Package	Functions
<code>readxl</code>	<code>read_excel('my-spreadsheet.xls', sheet = 1), read_xls('my-spreadsheet.xls'), read_xlsx('my-spreadsheet.xlsx')</code>
<code>readstata13</code>	<code>read.dta13('my-stata-data.dta')</code>
<code>readr</code>	<code>read_csv('my-csv-file.csv'), read_csv2('my-csv-file.csv'), read_delim(), read_rds()</code>
<code>vroom</code>	<code>vroom('my-csv-file.csv')</code>
<code>tidyxl</code>	<code>xlsx_cells('my_nightmare_file.xlsx')</code>
<code>haven</code>	<code>read_dta(), read_sas(), read_sav(), read_spss(), read_stata()</code>
<code>utils</code>	<code>read.csv, read.delim, read.table</code>

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 install.packages("haven")
2 install.packages("readstata13")
3 install.packages("tidyxl")
4 install.packages("readxl")
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 install.packages("haven")
2 install.packages("readstata13")
3 install.packages("tidyxl")
4 install.packages("readxl")
5
6 library (haven)
7 library (readstata13)
8 library (tidyxl)
9 library (readxl)
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 install.packages("haven")
2 install.packages("readstata13")
3 install.packages("tidyxl")
4 install.packages("readxl")
5
6 library (haven)
7 library (readstata13)
8 library (tidyxl)
9 library (readxl)
10
11 ?read_dta
12 ?read_xls
13 ?read.dta13
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library(readxl)
2 library(dplyr)
3
4 wvs_data <- read_xlsx("wvs_greatBritain.xlsx")
5 glimpse(wvs_data)
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
Rows: 2,609
Columns: 368
$ `version: Version of Data File`  

<chr> ...
$ `doi: Digital Object Identifier`  

<chr> ...
$ `A_YEAR: Year of survey`  

<chr> ...
$ `B_COUNTRY: ISO 3166-1 numeric country code`  

<chr> ...
$ `B_COUNTRY_ALPHA: ISO 3166-1 alpha-3 country code`  

<chr> ...
$ `C_COW_NUM: CoW country code numeric`  

<chr> ...
$ `C_COW_ALPHA: CoW country code alpha`  

<chr> ...
$ `D_INTERVIEW: Interview ID`  

<chr> ...
$ `J_INTPDATE: Date of interview`  

<chr> ...
$ `FW_START: Year/month of start-fieldwork`  

<chr> ...
$ `FW_END: Year/month of end-fieldwork`  

<chr> ...
```

*Numeric*  
*Integer*  
*Character*  
*Factor*  
*Logical*

Assess the structure of the data with `glimpse()` from `{dplyr}` 

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library(readxl)
2 library(dplyr)
3
4 wvs_data <- read_xlsx("wvs_greatBritain.xlsx")
5 head(wvs_data, 10)
```

- Remove duplicate or irrelevant observations

A tibble: 10 × 368			
Q38: It is children duty to take care of ill parent <chr>	Q39: People who don't work turn lazy <chr>	Q40: Work is a duty towards society <chr>	▶
Disagree	Strongly disagree	Strongly disagree	Strongly disagree
Disagree	Neither agree or disagree	Neither agree or disagree	Agree
Agree	Disagree	Agree	Agree
Disagree	Disagree	Disagree	Disagree
Neither agree nor disagree	Neither agree or disagree	Agree	Agree
Disagree	Disagree	Agree	Agree
Disagree	Neither agree or disagree	Don't know	Agree
Disagree	Disagree	Disagree	Disagree
Disagree	Disagree	Disagree	Neither agree or disagree
Disagree strongly	Disagree		

1-10 of 10 rows | 73-75 of 368 columns

- Handle (remove) missing data

- Validate

Assess the structure of the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate



# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Inspect the first 10 rows in the data with `head()` from `{utils}` 📁

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 library (readxl)
2 library(dplyr)
3 library (janitor)
4
5 wvs_data <- read_xlsx("wvs_greatBritain.xlsx")
6 head(wvs_data, 10)
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Fix column names with `clean_names()`  
from `{janitor}` 🗂️

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 library (readxl)
2 library(dplyr)
3 library (janitor)
4
5 wvs_data <- read_xlsx("wvs_greatBritain.xlsx")
6 head(wvs_data, 10)
7
8 ?clean_names
9
10 wvs_clean_data <- clean_names(wvs_data)
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Assess the structure of the data with  
`glimpse()` from `{dplyr}` 📁

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 library (readxl)
2 library(dplyr)
3 library (janitor)
4
5 wvs_data <- read_xlsx("wvs_greatBritain.xlsx")
6 head(wvs_data, 10)
7
8 ?clean_names
9
10 wvs_clean_data <- clean_names(wvs_data)
11 glimpse(wvs_clean_data)
```

### Cleaned column names

```
Rows: 2,609
Columns: 368
$ version_version_of_data_file
<chr> ...
$ doi_digital_object_identifier
<chr> ...
$ a_year_year_of_survey
<chr> ...
$ b_country_iso_3166_1_numeric_country_code
<chr> ...
$ b_country_alpha_iso_3166_1_alpha_3_country_code
<chr> ...
$ c_cow_num_co_w_country_code_numeric
<chr> ...
$ c_cow_alpha_co_w_country_code_alpha
<chr> ...
$ d_interview_interview_id
<chr> ...
$ j_intdate_date_of_interview
<chr> ...
$ fw_start_year_month_of_start_fieldwork
<chr> ...
```

### Uncleaned column names

```
Rows: 2,609
Columns: 368
$ `version: Version of Data File`
<chr> ...
$ `doi: Digital Object Identifier`
<chr> ...
$ `A_YEAR: Year of survey`
<chr> ...
$ `B_COUNTRY: ISO 3166-1 numeric country code`
<chr> ...
$ `B_COUNTRY_ALPHA: ISO 3166-1 alpha-3 country code`
<chr> ...
$ `C_COW_NUM: Cow country code numeric`
<chr> ...
$ `C_COW_ALPHA: Cow country code alpha`
<chr> ...
$ `D_INTERVIEW: Interview ID`
<chr> ...
$ `J_INTDATE: Date of interview`
<chr> ...
$ `FW_START: Year/month of start-fieldwork`
<chr> ...
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

*Keep only the relevant columns with  
select() from {dplyr}* 📁

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 library (readxl)
2 library(dplyr)
3 library (janitor)
4
5 wvs_data <- read_xlsx("wvs_greatBritain.xlsx")
6 head(wvs_data, 10)
7
8 ?clean_names
9
10 wvs_clean_data <- clean_names(wvs_data)
11
12 sub_wvs_data <- wvs_clean_data %>%
13   ## Select a few columns
14   select(a_year_year_of_survey,
15         q261_year_of_birth,
16         q260_sex,
17         h_urbrural_urban_rural,
18         q269_respondent_citizen,
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Rename columns with long or complicated names with `rename()` from `{dplyr}` 

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 sub_wvs_data <- wvs_clean_data %>%
2   ## Select a few columns
3   select(a_year_year_of_survey,
4         q261_year_of_birth,
5         q260_sex,
6         h_urbrural_urban_rural,
7         q269_respondent_citizen,
8         q223_local_party_preference_local_name,
9         q165_believe_in_god,
10        q191_justifiable_violence_against_other_people,
11        q275_highest_educational_level_respondent_isced_2011) %>%
12      ## Rename columns
13      rename (survey_yr = a_year_year_of_survey,
14               party_pref = q223_local_party_preference_local_name,
15               violence_just = q191_justifiable_violence_against_other_people,
16               education = q275_highest_educational_level_respondent_isced_2011,
17               residence = h_urbrural_urban_rural)
```

```
Rows: 2,609
Columns: 9
$ survey_yr          <chr> "2022", "2022", "2022", "2022", "2022", "2022", "2022", "2022", "2022", ...
$ q261_year_of_birth <chr> "1967", "1980", "2000", "1950", "1952", "1971", "1988", "1966", "1947", "1944", ...
$ q260_sex            <chr> "Female", "Female", "Female", "Male", "Female", "Female", "Female", "Female", ...
$ residence           <chr> "Rural", "Rural", "Rural", "Rural", "Urban", "Urban", "Urban", "Urban", "Urban", ...
$ q269_respondent_citizen <chr> "Yes", ...
$ party_pref          <chr> "4", "GBR: Labour Party", "GBR: Labour Party", "GBR: Liberal Democrats", ...
$ q165_believe_in_god <chr> "Yes", "Don't know", "No", "No", "Yes", "Yes", "Yes", "Yes", "Yes", "Yes", ...
$ violence_just        <chr> "Never justifiable", "Never justifiable", "Never justifiable", "Never justifiable", ...
$ education            <chr> "Upper secondary education (ISCED 3)", "Master or equivalent (ISCED 7)", "Post-secondary non-tertiary education ...
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Tabulate a few columns to understand the structure and identify potential structural errors with `table()` from `{base}` 📈

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 table (sub_wvs_data$q260_sex)
```

Female	
	1471
Male	
	1105
No answer	
	17
Other missing; Multiple answers Mail (EVS)	
	16

```
1 table (sub_wvs_data$party_pref)
```

-1	-2
261	111
-5	4
19	271
5	GBR: British National Party
19	3
GBR: Conservative and Unionist Party	GBR: Democratic Unionist Party
552	2
GBR: Green Party	GBR: Independence Party
178	20
GBR: Labour Party	GBR: Liberal Democrats
702	229
GBR: Plaid Cymru	GBR: Reform UK
40	24
GBR: Scottish National Party	GBR: Sinn Féin
175	1
---	---

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Create a new object `wvs_clean_1` from  
`wvs_clean_1`

```
1 wvs_clean_1 <- sub_wvs_data
```



- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Create a variable (replace if existing) with  
mutate() from {dplyr} 📂

```
1 wvs_clean_1 <- sub_wvs_data %>%
2   mutate (q260_sex = if_else((q260_sex != "Female" & q260_sex != "Male"),
3                               true = NA,
4                               false = q260_sex,
5                               missing = NA))
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

*Replace values that don't meet a condition to NA with `ifelse()` from {dplyr}* 📁

```
1 wvs_clean_1 <- sub_wvs_data %>%
2   mutate (q260_sex = if_else((q260_sex != "Female" & q260_sex != "Male"),
3                               true = NA,
4                               false = q260_sex,
5                               missing = NA))
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Replace all less meaningful values to missing (NA) with `ifelse()` from `{dplyr}` 

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wvs_clean_1 <- sub_wvs_data %>%
2   mutate (q260_sex = if_else((q260_sex != "Female" & q260_sex != "Male"),
3                             true = NA,
4                             false = q260_sex,
5                             missing = NA),
6
7   residence = if_else((residence == "No answer; Missing"),
8                        true = NA,
9                        false = residence,
10                       missing = NA),
11
12  q269_respondent_citizen = if_else((q269_respondent_citizen != "No" &
13                                     q269_respondent_citizen != "Yes"),
14                                     true = NA,
15                                     false = q269_respondent_citizen,
16                                     missing = NA),
17
18  q261_year_of_birth = if_else((q261_year_of_birth == "No answer" |
19                               q261_year_of_birth == "Other missing",
20                               true = NA,
21                               false = q261_year_of_birth).
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Assess the structure of the data, again with `glimpse()` and tabulate with `table()`

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 glimpse(wvs_clean_1$q261_year_of_birth)
chr [1:2609] "1967" "1980" "2000" "1950" "1952" "1971" "1988" "1966" ...
```

```
1 table(wvs_clean_1$violence_just)
```

2	3	4	5
277	168	63	107
6	7	8	Always justifiable
22	21	5	12
Never justifiable			
1913			

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Create a variable (replace if existing) with  
`mutate()` from `{dplyr}` 📂

- Remove duplicate or irrelevant observations

```
1 wvs_clean_2 <- wvs_clean_1 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence))
```

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Convert values in a variable to numeric with `as.numeric()` and to categories with `as.factor()` from `{base}` 📂

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wvs_clean_2 <- wvs_clean_1 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence))
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Recode values in a variable with  
`case_when()` from `{dplyr}` 

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wvs_clean_3 <- wvs_clean_2 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence)) %>%
7   mutate (violence_just = case_when(violence_just == "Always justifiable" ~ 9,
8                                     violence_just == "Never justifiable" ~ 1,
9                                     .default = as.numeric(violence_just)))
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Assess the structure of the data with  
`glimpse()` from `{dplyr}` 📁

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wvs_clean_3 <- wvs_clean_2 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence)) %>%
7   mutate (violence_just = case_when(violence_just == "Always justifiable" ~ 9,
8                                     violence_just == "Never justifiable" ~ 1,
9                                     .default = as.numeric(violence_just)))
```

num [1:2609] 1 1 1 1 1 1 1 1 1 1 ...

```
1 table (wvs_clean_3$violence_just)
```

1	2	3	4	5	6	7	8	9
1913	277	168	63	107	22	21	5	12

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wvs_clean_3 <- wvs_clean_2 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence)) %>%
7   mutate (violence_just = case_when(violence_just == "Always justifiable" ~ 9,
8                                     violence_just == "Never justifiable" ~ 1,
9                                     .default = as.numeric(violence_just)))
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Create a variable age (replace if existing)  
with `mutate()` from `{dplyr}` 📂

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wvs_clean_3 <- wvs_clean_2 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence)) %>%
7   mutate (violence_just = case_when(violence_just == "Always justifiable" ~ 9,
8                                     violence_just == "Never justifiable" ~ 1,
9                                     .default = as.numeric(violence_just))) %>%
10
11   mutate (age = survey_yr - q261_year_of_birth)
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

Assess the distribution of age with  
`table()` from `{base}` 📈

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

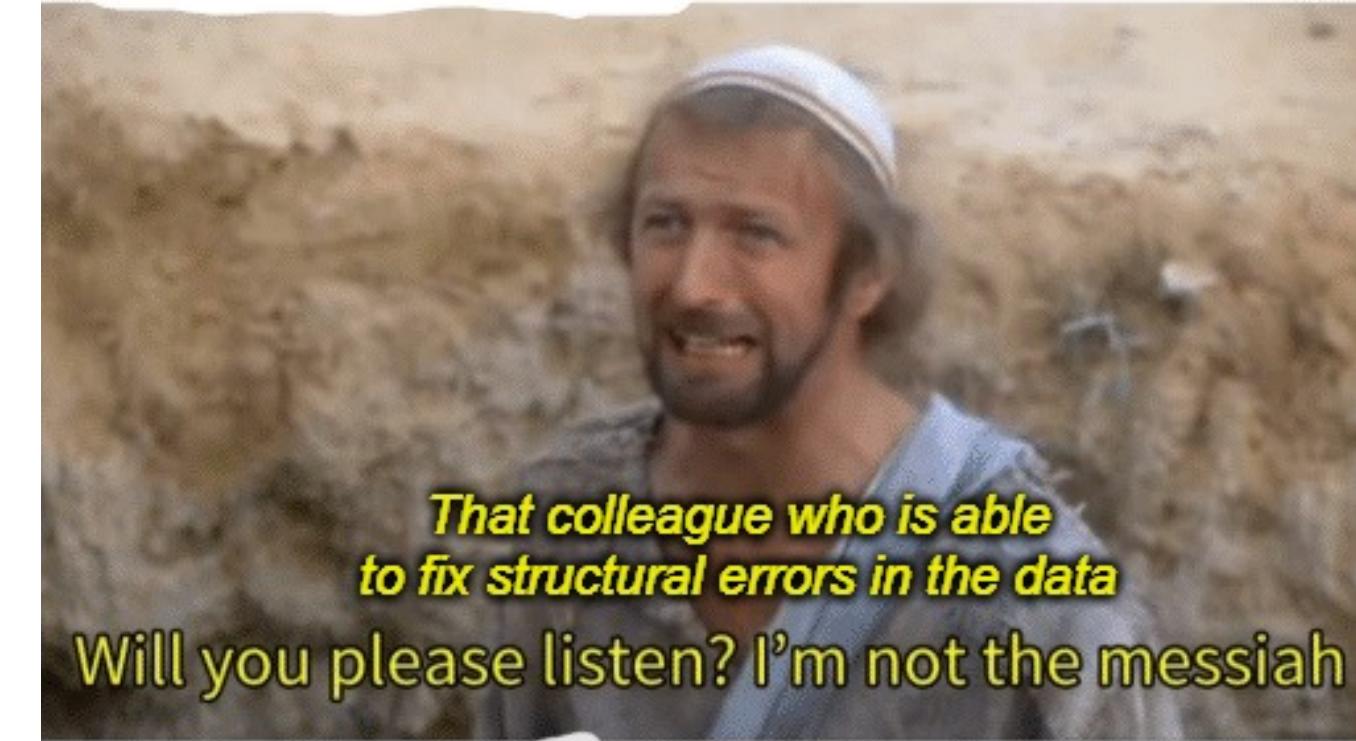
```
1 wvs_clean_3 <- wvs_clean_2 %>%
2   mutate (q261_year_of_birth = as.numeric(q261_year_of_birth),
3           survey_yr = as.numeric(survey_yr),
4           q260_sex = as.factor(q260_sex),
5           q269_respondent_citizen = as.factor(q269_respondent_citizen),
6           residence = as.factor(residence)) %>%
7   mutate (violence_just = case_when(violence_just == "Always justifiable" ~ 9,
8                                     violence_just == "Never justifiable" ~ 1,
9                                     .default = as.numeric(violence_just))) %>%
10
11   mutate (age = survey_yr - q261_year_of_birth)
12
13 table (wvs_clean_3$age)
```

```
19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44
25 15 16 30 33 33 22 33 36 27 29 40 48 58 41 45 36 32 45 37 43 35 46 33 45 42
45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70
39 43 29 40 30 34 35 43 38 41 47 40 46 39 43 52 37 43 22 45 47 49 34 42 37 47
71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 94 95 96 99
42 31 42 39 39 38 22 40 18 27 19 13 17 19 18 15 5 10 7 6 7 3 1 3 3 1
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate



imgflip.com

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- *Remove duplicate or irrelevant observations*
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 wvs_clean_4 <- wvs_clean_3 %>%
2   filter (age >= 18 & age <= 34)
3
4 str (wvs_clean_4)

tibble [531 × 10] (S3:tbl_df/tbl/data.frame)
$ survey_yr           : num [1:531] 2022 2022 2022 2022 2022 ...
$ q261_year_of_birth  : num [1:531] 2000 1988 1993 1996 1990 ...
$ q260_sex             : Factor w/ 2 levels "Female","Male": 1 1 2 2 2 1 1 1 1 ...
$ residence            : Factor w/ 2 levels "Rural","Urban": 1 2 2 2 2 1 1 1 2 ...
$ q269_respondent_citizen: Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 2 2 2 ...
$ party_pref           : chr [1:531] "GBR: Labour Party" NA "GBR: Labour Party" "GBR: Conservative
and Unionist Party" ...
$ q165_believe_in_god  : chr [1:531] "No" "Yes" "No" "No" ...
$ violence_just         : num [1:531] 1 1 5 3 6 1 1 1 3 3 ...
$ education             : chr [1:531] "Post-secondary non-tertiary education (ISCED 4)" "Upper
secondary education (ISCED 3)" "Bachelor or equivalent (ISCED 6)" "Bachelor or equivalent (ISCED 6)"
...
tibble [2,609 × 10] (S3:tbl_df/tbl/data.frame)
$ survey_yr           : num [1:2609] 2022 2022 2022 2022 2022 ...
$ q261_year_of_birth  : num [1:2609] 1967 1980 2000 1950 1952 ...
$ q260_sex             : Factor w/ 2 levels "Female","Male": 1 1 1 1 2 1 1 1 1 ...
$ residence            : Factor w/ 3 levels "No answer; Missing",...: 2 2 2 2 2 3 3 3 3 ...
$ q269_respondent_citizen: Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 2 2 2 ...
$ party_pref           : chr [1:2609] NA "GBR: Labour Party" "GBR: Labour Party" "GBR: Liberal Democrats" ...
$ q165_believe_in_god  : chr [1:2609] "Yes" NA "No" "No" ...
$ violence_just         : num [1:2609] 1 1 1 1 1 1 1 1 ...
$ education             : chr [1:2609] "Yes" "Yes" "Yes" "Yes" ...
$ age                  : num [1:2609] 55 42 22 72 70 51 34 56 75 78 ...
```

- Remove duplicate or irrelevant observations
- Keep only data (or observations) from young adults aged 18-34 years with `filter()` from `{dplyr}` 🗂
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- Remove duplicate or irrelevant observations
- *Handle (remove) unwanted outliers*
- Handle (remove) missing data
- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 table(wvs_clean_4$education)
```

Bachelor or equivalent (ISCED 6)	150
Doctoral or equivalent (ISCED 8)	4
Early childhood education (ISCED 0) / no education	4
Lower secondary education (ISCED 2)	105
Master or equivalent (ISCED 7)	75
Post-secondary non-tertiary education (ISCED 4)	11
Primary education (ISCED 1)	3
Short-cycle tertiary education (ISCED 5)	50

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

*mmmCheck if there are any missing values in the data with `anyNA()` from `{base}`*



- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

mmmCheck if there are any missing values in the data with `anyNA()` from `{base}`



- Handle (remove) missing data

- Validate

```
1 wvs_clean_4 <- wvs_clean_4 %>%
2   mutate (education = case_when((education == "Early childhood education (ISCED 0) /
3
4   education == "Primary education (ISCED 1)" |
5   education == "Upper secondary education (ISCED 3)" |
6   education == "Lower secondary education (ISCED 2)" |
7
8   education == "Short-cycle tertiary education (ISCED 4)" |
9   education == "Post-secondary non-tertiary education (ISCED 5)" |
10  education == "Bachelor or equivalent (ISCED 6)" |
11  education == "Master or equivalent (ISCED 7)" |
12  education == "Doctoral or equivalent (ISCED 8)") |
13
14  .default = factor(education)) )
15
16 table (wvs_clean_4$education)
```

	Post-secondary	Secondary or less
	290	221

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 table(wvs_clean_4$violence_just)
```

1	2	3	4	5	6	7	8	9
330	69	51	25	32	5	9	2	2

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

*mmmCheck if there are any missing values in the data with `anyNA()` from `{base}`*



- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 wvs_clean_4 <- wvs_clean_4 %>%
2   mutate (violence_just = if_else((violence_just >= 2),
3   true = "Justified",
4   false = "Never justified",
5   missing = NA)) %>%
6   mutate (violence_just = as.factor(violence_just))
7
8 table (wvs_clean_4$violence_just)
```

	Justified	Never justified
	195	330

- Handle (remove) unwanted outliers

*mmmCheck if there are any missing values in the data with `anyNA()` from `{base}`*



- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- *Handle (remove) missing data*
- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 anyNA(wvs_clean_4)
```

```
[1] TRUE
```

- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

*Check if there are any missing values in the data with `anyNA()` from `{base}`* 📁

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library (visdat)
2 vis_miss(wvs_clean_4)
```

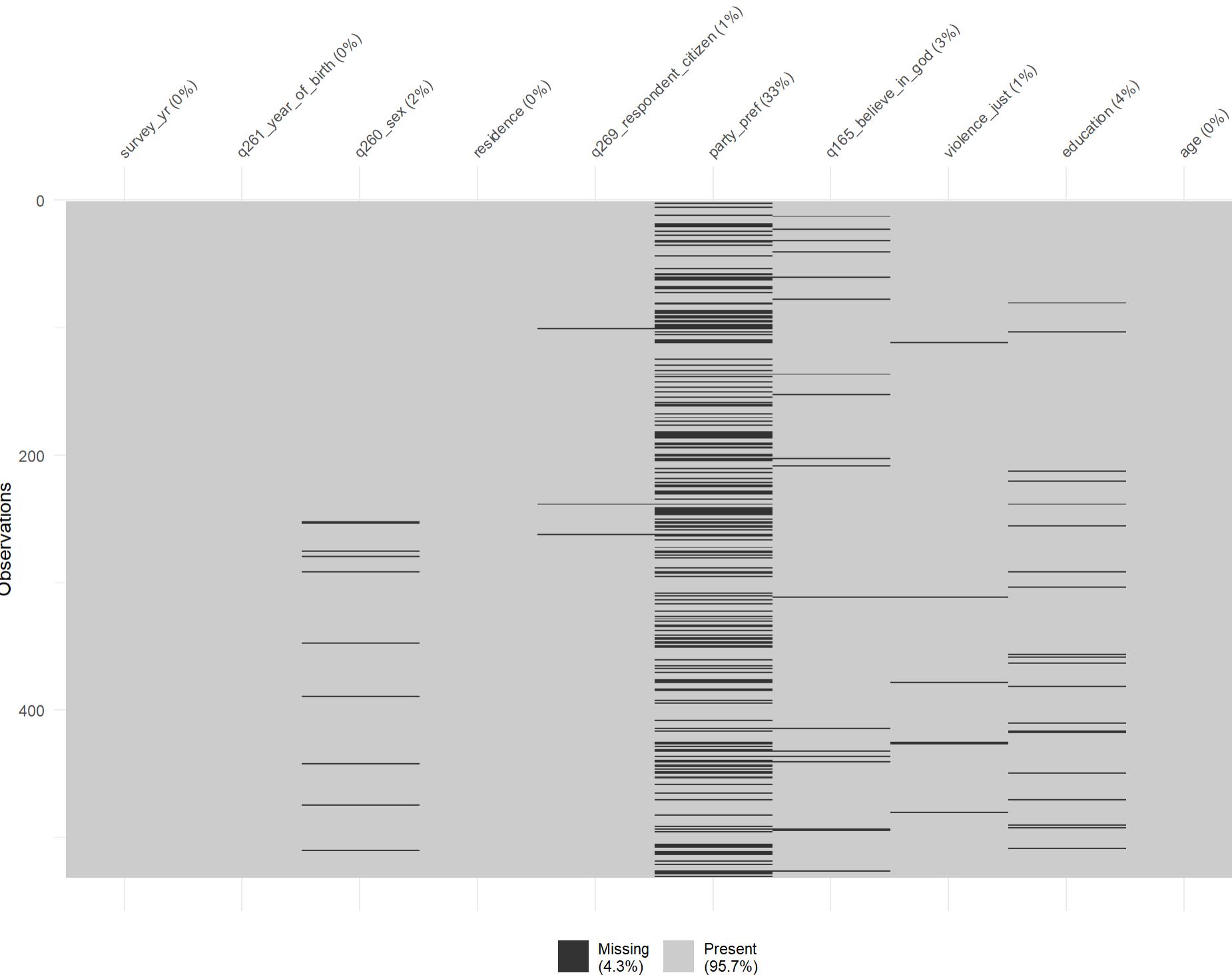
- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

Visualise missing values across the dataset  
with `vis_miss()` from `{visdat}` 🗂️

- Validate



# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library (visdat)
2 vis_miss(wvs_clean_4)
3
4 wvs_clean_5 <- wvs_clean_4 %>%
5   filter (!is.na (q260_sex) &
6         !is.na(residence) &
7         !is.na(q269_respondent_citizen))
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

*Remove missing values in each column with  
filter() from {dplyr}* 🗂

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library (visdat)
2 vis_miss(wvs_clean_4)
3
4 wvs_clean_5 <- wvs_clean_4 %>%
5   filter (!is.na (q260_sex) &
6         !is.na(residence) &
7         !is.na(q269_respondent_citizen))
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

Check whether a column has missing values with `is.na()` from `{base}` 📈

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library (visdat)
2 vis_miss(wvs_clean_5)
```

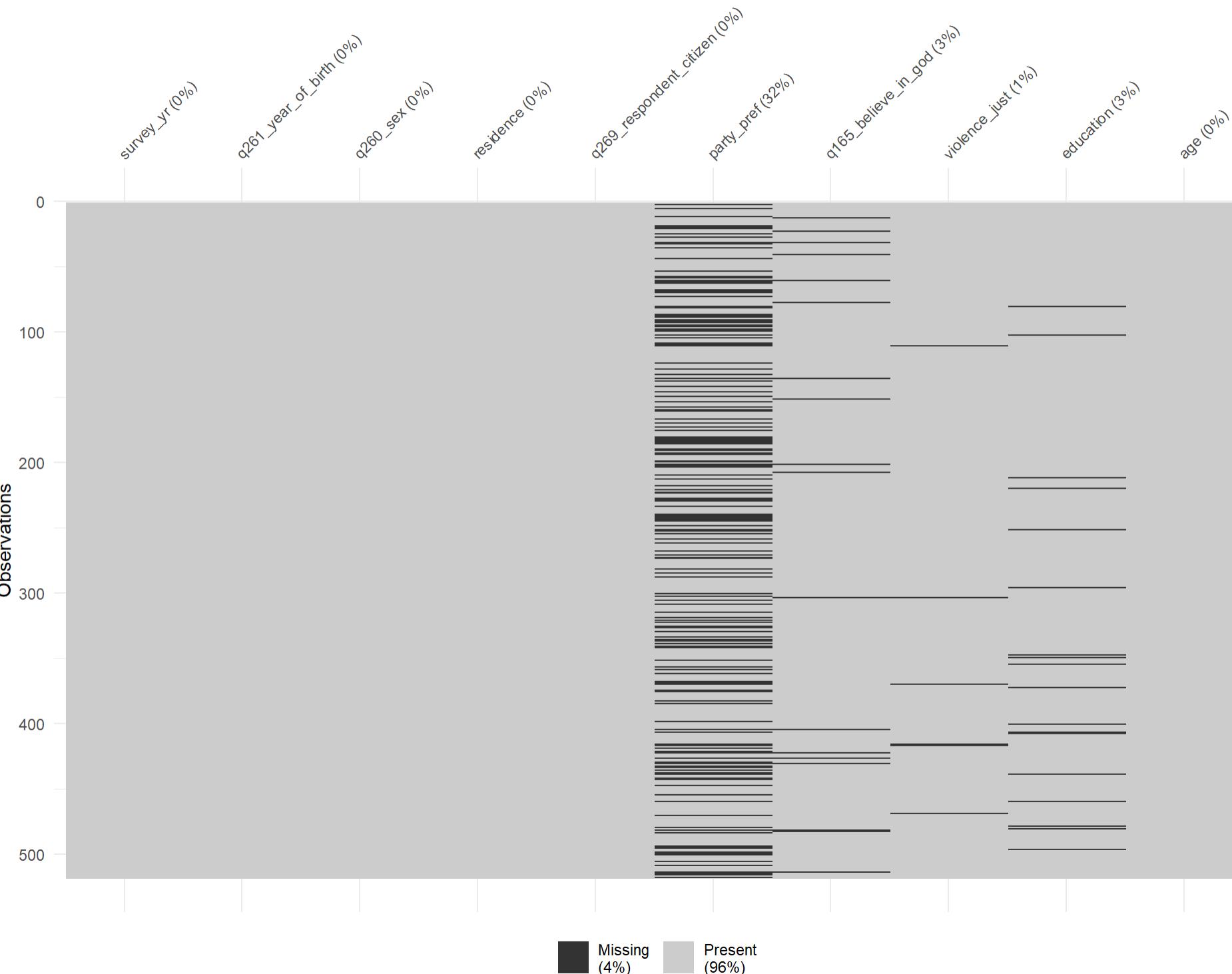
- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

Visualise missing values across the dataset  
with `vis_miss()` from `{visdat}` 🗂️

- Validate



# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library (visdat)
2 vis_miss(wvs_clean_4)
3
4 wvs_clean_5 <- wvs_clean_4 %>%
5   filter (!is.na (q260_sex) &
6         !is.na(residence) &
7         !is.na(q269_respondent_citizen))
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

*Remove missing values in each column with  
filter() from {dplyr}* 🗂

- Validate

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

*Remove missing values in each column with  
`filter()` from {dplyr} 📂*

- Validate

```
1 library (visdat)
2 vis_miss(wvs_clean_4)
3
4 wvs_clean_5 <- wvs_clean_4 %>%
5   filter (!is.na (q260_sex) &
6         !is.na(residence) &
7         !is.na(q269_respondent_citizen)) %>%
8   filter (!is.na (party_pref) &
9         !is.na(q165_believe_in_god) &
10        !is.na(violence_just) &
11        !is.na(education))
```

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 library (visdat)
2 vis_miss(wvs_clean_5)
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- **Handle (remove) missing data**

Visualise missing values across the dataset  
with `vis_miss()` from `{visdat}` 🗂️

- Validate



# {Tidy} data management:

## A case study using the [World Value Survey](#)

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- *Validate*



# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 ## Sample of young people with complete cases  
2 dim(wvs_clean_5)
```

[1] 332 10

```
1 ## Full sample of young people with missing cases  
2 dim(wvs_clean_4)
```

[1] 531 10

```
1 ## Full adult sample in the dataset with missing cases  
2 dim(wvs_clean_3)
```

[1] 2609 10

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- **Validate**

Assess the data dimensions with `dim()`  
from `{base}` 📂

# {Tidy} data management:

## A case study using the World Value Survey

- Fix structural errors

```
1 ## Sample of young people with complete cases  
2 dim(wvs_clean_5)
```

[1] 332 10

```
1 ## Full sample of young people with missing cases  
2 dim(wvs_clean_4)
```

[1] 531 10

```
1 ## Full adult sample in the dataset with missing cases  
2 dim(wvs_clean_3)
```

[1] 2609 10

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- **Validate**

Assess the data dimensions with `dim()`  
from `{base}` 📂

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 ## Sample of young people with complete cases  
2 dim(wvs_clean_5)
```

[1] 332 10

```
1 ## Full sample of young people with missing cases  
2 dim(wvs_clean_4)
```

[1] 531 10

```
1 ## Full adult sample in the dataset with missing cases  
2 dim(wvs_clean_3)
```

[1] 2609 10

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- **Validate**

Assess the data dimensions with `dim()`  
from `{base}` 📂

# {Tidy} data management

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

The screenshot shows the World Bank DataBank interface for the Environment Social and Governance (ESG) data. The top navigation bar includes the World Bank logo, a home icon, a sign-in button, and links for 'Metadata' and 'Download' (Excel, CSV, Tabbed TXT). The main title is 'DataBank | Environment Social and Governance'. Below the title, there are tabs for 'Table' (selected), 'Chart', and 'Map'. On the left, there's a sidebar with 'Variables' (Available 85, Selected 1), 'Database' (Available 239, Selected 239), and 'Series' (Available 71, Selected 71). A search bar says 'Enter Keywords for' with a magnifying glass icon. A dropdown menu lists variables: 'Access to clean fuels and technologies for cooking (% of population)', 'Access to electricity (% of population)', 'Adjusted savings: natural resources depletion (% of GNI)', 'Adjusted savings: net forest depletion (% of GNI)', 'Agricultural land (% of land area)', and 'Agriculture, forestry, and fishing, value added (% of GDP)'. The value for Agricultural land is 58.7. At the bottom right is a 'Help/Feedback' button.

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 library(dplyr)
2 wb_dt <- read.csv("wb_databank.csv")
3
4 head(wb_dt, 10)
5 tail(wb_dt, 10)
```

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

	Country.Name	Country.Code	Series.Name	Series
16964	World	WLD	Strength of legal rights index (0=weak to 12=strong)	IC.LGL
16965	World	WLD	Terrestrial and marine protected areas (% of total territorial ...	ER.PTI
16966	World	WLD	Tree Cover Loss (hectares)	AG.LN
16967	World	WLD	Unemployment, total (% of total labor force) (modeled ILO ...	SL.UEI
16968	World	WLD	Unmet need for contraception (% of married women ages 1...	SP.UM
16969	World	WLD	Voice and Accountability: Estimate	VA.ES
16970				
16971				
16972				
16973	Data from database: Environment Social and Governance (E...			
16974	Last Updated: 10/02/2023			
16975	Code	License Type	Indicator Name	Short
16976	EG.CFT.ACCTS.ZS	CC BY-4.0	Access to clean fuels and technologies for cooking (% of po...	
16977	EG.ELC.ACCTS.ZS	CC BY-4.0	Access to electricity (% of population)	
16978	NY.ADJ.DRES.GN.ZS	CC BY-4.0	Adjusted savings: natural resources depletion (% of GNI)	
16979	NY.ADJ.DFOR.GN.ZS	CC BY-4.0	Adjusted savings: net forest depletion (% of GNI)	
16980	AG.LND.AGRI.ZS	CC BY-4.0	Agricultural land (% of land area)	
16981	NV.AGR.TOTL.ZS	CC BY-4.0	Agriculture, forestry, and fishing, value added (% of GDP)	
16982	ER.H2O.FWTL.ZS	CC BY-4.0	Annual freshwater withdrawals, total (% of internal resources)	
16983	SI.SPR.PCAP.ZG	CC BY-4.0	Annualized average growth rate in per capita real survey me...	The gi
16984	SH.DTH.COMM.ZS	CC BY-4.0	Cause of death, by communicable diseases and maternal, pr...	
16985	SL.TLF.0714.ZS	CC BY-4.0	Children in employment, total (% of children ages 7-14)	
16986				

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 library(dplyr)
2
3 wb_dt <- read.csv("wb_databank.csv")
4
5 head(wb_dt, 10)
6
7 tail(wb_dt, 10)
8
9 wb_data <- wb_dt[1:16969,]
```

Inspect the first 10 rows in the data with `head()` from `{utils}` 

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 wb_data <- wb_dt[1:16969, ]  
2  
3 wb_data_2 <- wb_data %>%  
4     janitor::clean_names()
```

Inspect the first 10 rows in the data with `head()` from `{utils}` 

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 wb_data <- wb_dt[1:16969, ]
2
3 wb_data_2 <- wb_data %>%
4   janitor::clean_names() %>%
5   pivot_longer(cols = contains("_yr"),
6                 values_to = "values",
7                 names_to = "year")
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wb_data <- wb_dt[1:16969, ]
2
3 wb_data_2 <- wb_data %>%
4   janitor::clean_names() %>%
5   pivot_longer(cols = contains("_yr"),
6                 values_to = "values",
7                 names_to = "year") %>%
8   mutate(values = ifelse(values == "...",
9                         NA, values)) %>%
10  mutate(values = as.numeric(values))
```

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wb_data <- wb_dt[1:16969, ]
2
3 wb_data_2 <- wb_data %>%
4   janitor::clean_names() %>%
5   pivot_longer(cols = contains("_yr"),
6                 values_to = "values",
7                 names_to = "year") %>%
8   mutate(values = ifelse(values == "...",
9                         NA, values)) %>%
10  mutate(values = as.numeric(values)) %>%
11  mutate(period = str_extract(year, "[0-9]+")) %>%
12  mutate(year = as.numeric(period))
```

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 new_dta <- wb_data_2
```

- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` 

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 new_dta <- wb_data_2 %>%
2   filter (year >= 2000 & year <= 2020) %>%
3   filter (!is.na (values))
```

- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` 

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 data_key <- new_dta %>%
2     select (series_name, series_code)
```

- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 data_key <- new_dta %>%
 2   select (series_name, series_code) %>%
 3   reframe(series_code = unique(series_code),
 4           series_name = unique(series_name))
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 data_key <- new_dta %>%
  select (series_name, series_code) %>%
  reframe(series_code = unique(series_code),
          series_name = unique(series_name)) %>%
  mutate (series_code = str_to_lower(series_code)) %>%
  mutate (series_code = str_replace_all(series_code, "\\\.",
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 wide_dta <- new_dta %>%
2     select (country_name, year,
3             series_code, values)
```

- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 wide_dta <- new_dta %>%
2   select (country_name, year,
3           series_code, values) %>%
4   pivot_wider(names_from = series_code,
5               values_from = values)
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` 

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 wide_dta <- new_dta %>%
2   select (country_name, year,
3           series_code, values) %>%
4   pivot_wider(names_from = series_code,
5               values_from = values) %>%
6   janitor::clean_names()
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Inspect the first 10 rows in the data with `head( )` from `{utils}` ↗

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

```
1 wide_dta <- new_dta %>%
2   select (country_name, year,
3           series_code, values) %>%
4   pivot_wider(names_from = series_code,
5               values_from = values) %>%
6   janitor::clean_names() %>%
7   select (country_name, year, eg_cft_accs_zs,
8           eg_elc_accs_zs, en_atm_co2e_pc,
9           en_clc_heat_xd, sp_dyn_tfrt_in )
```

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

Inspect the first 10 rows in the data with `head()` from `{utils}` 

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wide_dta <- new_dta %>%
2   select (country_name, year,
3           series_code, values) %>%
4   pivot_wider(names_from = series_code,
5               values_from = values) %>%
6   janitor::clean_names() %>%
7   select (country_name, year, eg_cft_accs_zs,
8           eg_elc_accs_zs, en_atm_co2e_pc,
9           en_clc_heat_xd, sp_dyn_tfrt_in) %>%
10  mutate (eg_cft_accs_zs = round (eg_cft_accs_zs, 2),
11          eg_elc_accs_zs = round (eg_elc_accs_zs, 2),
12          en_atm_co2e_pc = round (en_atm_co2e_pc, 2),
13          en_clc_heat_xd = round (en_clc_heat_xd, 2),
14          sp_dyn_tfrt_in = round (sp_dyn_tfrt_in, 2))
```

Inspect the first 10 rows in the data with `head()` from `{utils}` 

# {Tidy} data management:

## A case study using the World Bank Data

- Fix structural errors

- Remove duplicate or irrelevant observations

- Handle (remove) unwanted outliers

- Handle (remove) missing data

- Validate

```
1 wide_dta <- new_dta %>%
2   select (country_name, year,
3           series_code, values) %>%
4   pivot_wider(names_from = series_code,
5               values_from = values) %>%
6   janitor::clean_names() %>%
7   select (country_name, year, eg_cft_accs_zs,
8           eg_elc_accs_zs, en_atm_co2e_pc,
9           en_clc_heat_xd, sp_dyn_tfrt_in) %>%
10  mutate (eg_cft_accs_zs = round (eg_cft_accs_zs, 2),
11          eg_elc_accs_zs = round (eg_elc_accs_zs, 2),
12          en_atm_co2e_pc = round (en_atm_co2e_pc, 2),
13          en_clc_heat_xd = round (en_clc_heat_xd, 2),
14          sp_dyn_tfrt_in = round (sp_dyn_tfrt_in, 2))
```

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

# {Tidy} data management:

## A case study using the **World Bank Data**

- Fix structural errors
- Remove duplicate or irrelevant observations
- Handle (remove) unwanted outliers
- Handle (remove) missing data
- Validate

```
1 final <- w:  
1   Tips for finding data  
2   | What type of data is needed?  
3   | Who could have collected the data?  
4  
5   Where to Find Data by:  
6   | Academics  
7   | Not-for-profit Organizations  
8   | Government Departments  
9   | Private/Commercial Companies
```

Inspect the first 10 rows in the data with `head()` from `{utils}` ↗

