# Tidying Data in R

Bas Machielsen

Utrecht University - R Cafe

December 28, 2019

# What is tidy data?

- In order for your data to be ready for analysis, it needs to be **tidy**.

- Tidy data is defined as:
  *data sets that are arranged such that each variable is a column and each observation (or case) is a row.*

- I will attempt to explain how to put this principle in practice using a set of packages called the **Tidyverse**, but more specifically, the `tidyr` package.

- You can download the .Rmd and .csv files from my github, and follow along!

# What is (un)tidy data?

*All tidy data are similar to each other. All untidy data are untidy in their own way.*

- Although tidy data shares the properties mentioned in the previous slide, untidy data refers (broadly speaking) to unstructured data.
- An example (from www.tidyverse.org):

Table 1: This is what an untidy dataset looks like

| name | treatmenta | treatmentb |
|------|-----------|-----------|
| John Smith | NA | 18 |
| Jane Doe | 4 | 1 |
| Mary Johnson | 6 | 7 |

# An example of untidy data

- Is each obervation in a row?

- Is each variable in a column?

- What are the variables here?

Table 2: This is also what an untidy dataset looks like

| treatment | John.Smith | Jane.Doe | Mary.Johnson |
|-----------|------------|----------|--------------|
| a | NA | 4 | 6 |
| b | 18 | 1 | 7 |

# Example: from untidy to tidy data

- The variables are (i) the treatment, (ii) the subject, and (iii) the 'score' corresponding to each subject-treatment observations.

```
preg %>%
    pivot_longer(treatmenta:treatmentb,
                 names_to = "treatment", values_to = "score") %>%
    mutate(treatment = gsub("treatment", "", treatment)) %>%
    arrange(name, treatment) %>%
    kable(caption = "Test", booktabs = TRUE,
          row.names = FALSE) %>%
    kable_styling(latex_options = "striped")
```

# Example: from untidy to tidy data

- The former chunk of code transforms the data in the first table to the following table:

Table 3: This is tidy data

| name | treatment | score |
|------|-----------|-------|
| Jane Doe | a | 4 |
| Jane Doe | b | 1 |
| John Smith | a | NA |
| John Smith | b | 18 |
| Mary Johnson | a | 6 |
| Mary Johnson | b | 7 |

# Transforming the data

- The key command in the former chunk of code is `pivot_longer`. With the arguements data, columns, names_to and values_to.

- All the other functions are essentially layout changes, not fundamental transformations.

- The command pivot_longer basically transforms the dataset from:

Table 4: This is what an untidy dataset looks like

| name | treatmenta | treatmentb |
|------|------------|------------|
| John Smith | NA | 18 |
| Jane Doe | 4 | 1 |
| Mary Johnson | 6 | 7 |

# Transforming the data

- to :

```r
pivot_longer(data = preg, cols = 2:3, names_to = "treatment",
             values_to = "score") %>%
   kable(caption = "This is what a tidy dataset looks like",
      booktabs = TRUE, row.names = FALSE) %>%
 kable_styling(latex_options = "striped")
```

Table 5: This is what a tidy dataset looks like

| name | treatment | score |
|------|-----------|-------|
| John Smith | treatmenta | NA |
| John Smith | treatmentb | 18 |
| Jane Doe | treatmenta | 4 |
| Jane Doe | treatmentb | 1 |
| Mary Johnson | treatmenta | 6 |
| Mary Johnson | treatmentb | 7 |

# Contents

- The best way to see how tidy data works is by using examples. In this lecture, I will demonstrate extensively how to create tidy data that is ready for analysis using two examples.

- First, I will show how to combine and tidy data on economic development from the World Bank.

- Afterwards, I show how to combine and tidy data from ORBIS, a commercial database.

- On the fly, I will demonstrate how simple it is to analyze or create graphs with tidy data.

- Finally, I will use other examples to show how to tidy more messy data, inspired by examples from the **tidyverse website**

# World Bank Dataset

- Let us now download some data from the World Bank database. I extract the 55 basic World Development Indicators from **here** for 20 countries.

- Importing the data..

```
worldbank <- read_csv("wb.csv")

dim(worldbank)
```

```
## [1] 1105   14
```

- The data has 1105 observations and 14 variables. That is not at all what we wanted. What do the data look like?

# World Bank Dataset

```r
kable(worldbank[1:4,c(1,2,4,5)],
      caption = "Work Bank Data - Untidy!",
      booktabs = TRUE, row.names = FALSE) %>%
  kable_styling(latex_options = "striped")
```

Table 6: Work Bank Data - Untidy!

| Country Name | Country Code | Series Code | 2009 [YR2009] |
|--------------|--------------|-------------|---------------|
| Argentina | ARG | SP.ADO.TFRT | 63.1896 |
| Argentina | ARG | NV.AGR.TOTL.ZS | 5.27362346890139 |
| Argentina | ARG | ER.H2O.FWTL.ZS | .. |
| Argentina | ARG | SH.STA.BRTC.ZS | 97.9 |

- Very untidy dataset! NA observations are entered as .., and variable names require an extensive definition before you know what they mean.

- Furthermore, years are not notated straightforwardly, and the data violates the tidy data principles.

- First, let us try to save the variable names and series code to another dataset (step 1), and delete the names from the worldbank dataset (step 2).

```
#Step 1
wbnames <- cbind(unique(worldbank[,3]),unique(worldbank[,4])) %>%
    na.omit()

#Step 2
worldbank <- worldbank[1:1100,-3]
```

# Almost there..

- Now, let's try to unpivot the columns containing the years..

```r
worldbank <- pivot_longer(worldbank, 4:13,
            names_to = "years",
            values_to = "value")
```

- This is what the data looks like right now.

Table 7: The data then looks like this.

| Country Code | Series Code | years | value |
|---|---|---|---|
| ARG | SP.ADO.TFRT | 2009 [YR2009] | 63.1896 |
| ARG | SP.ADO.TFRT | 2010 [YR2010] | 63.3154 |
| ARG | SP.ADO.TFRT | 2011 [YR2011] | 63.4412 |
| ARG | SP.ADO.TFRT | 2012 [YR2012] | 63.567 |

# Final steps

- Let's now clean the "years" variable (step 1), set the .. observations to NA (step 2), and "widen" the data so that every separate indicator gets a column (step 3).

- Step 3 is done using the pivot_wider command: it widens the data and shortens the number of observations, because it transfers information from rows to columns.

```r
#Step 1
worldbank <- worldbank %>%
    mutate(years = substring(years, 1, 4))

#Step 2
worldbank$value[worldbank$value == ".."] <- NA

#Step 3
worldbank <- pivot_wider(data = worldbank,
                         names_from = `Series Code`,
                         values_from = `value`)
```

## Done!

- `pivot_wider`, takes three arguments: the data, from which column it should take the new column names (names_from), and from which column it should take the values belonging to the corresponding cels (values_from).

- This is what it looks like!

Table 8: Tidy Data

| Country Name | Country Code | years | SP.ADO.TFRT | NV.AGR.TOTL.ZS |
|---|---|---|---|---|
| Argentina | ARG | 2009 | 63.1896 | 5.27362346890139 |
| Argentina | ARG | 2010 | 63.3154 | 7.13216745078964 |
| Argentina | ARG | 2011 | 63.4412 | 6.99873377022519 |
| Argentina | ARG | 2012 | 63.567 | 5.7817442068501 |

# Some analysis

- Suppose I want to know the military expenditures as a percentage of GDP MS.MIL.XPND.GD.ZS for each country in the dataset in 2015.

- I can use the `wbnames` information we've saved to add labels to the data.

```
worldbank %>%
    filter(years == "2015") %>%
    ggplot(aes(x = `Country Code`, y = `MS.MIL.XPND.GD.ZS`)) +
    geom_col() +
    labs(y = as.character(wbnames[match("MS.MIL.XPND.GD.ZS",
                                        wbnames[,2]),1])) +
    theme(axis.text.y =element_blank())
```