

# Week 3: 'Tidying' data

---

Joshua Rosenberg, Alex Lishinski

February 4, 2021

# Welcome!

---

Welcome to *week 3*!

Great job last week.

Now we feel like we have enough of an introduction to some of the basics of working with R that we can move to the more substantial material. However, don't be discouraged if you don't feel rock solid about any of the basics we've covered so far. We will continue to reinforce these as we move forward.

**Record the meeting**

# Course organization

---

- Website
  - Syllabus
  - Presentations
  - Homework
  - Videos (recordings of class and about specific topics)
- Slack: <https://datascienceinedu.slack.com/home>
- Zoom: <https://tennessee.zoom.us/j/98262405076>
- Email: jrosenb8@utk.edu and alishins@utk.edu

# Breakout rooms!

---

Starting with whomever has the most pets (tiebreaker: aggregate pet weight)...

- From the chapter introducing educational data science, which of the educational data science roles do you see as the most similar to the work you have done, are doing, or intend to do, and why? You can consult the chapter if necessary. If you are coming from outside of education, discuss any similarities (or the lack thereof)!
- What are you looking forward to (general life question)?

(10 minutes)

# Foundational R skills

---

A general framework for you to use as a foundation and as a set of concepts to help you work through the class.

The four core concepts we will use to build our framework are:

1. Projects
2. Packages
3. Data
4. Functions

You will use each of these in most of your analyses with R.

# Foundational R skills

---

These foundational skills are what we will build on as we move to data wrangling and tidying

We're not done building your understanding of working with packages, data, and functions

**Course schedule note, we're shifting the calendar around to move data cleaning and wrangling ahead of data visualization**

# Packages quick review

---

## 1. What are packages?

- Code bundles that add functionality to R
- Examples: ggplot2, dplyr, rtweet, quanteda, lme4

## 2. Where do we get packages?

- CRAN or GitHub

## 3. How do we install packages?

- `install.packages("pkg-name")`

## 4. How do we know what packages to use?

- Searching
- People and news related to R (more later – there are *tons*)
- CRAN task views

## 5. How do we use packages?

- `library(pkgname)`

# Overview of today's class

## Wrangling and Tidying

1. Reading in Data
2. Tidying data
3. Our Tidy data tools



# 1. Reading in Data

---

So far, we have used *built-in data*. There is a lot of built-in data!

Loading different types of data, assigning to a name

Comma-separated values (**.csv**)

```
library(readr)

data <- readr::read_csv(here("data", "answer_export.csv"))
```

# 1. Reading in Data

---

`.xlsx`

```
library(readxl)  
  
data <- read_excel((here("data", "schedule.xlsx")))
```

# 1. Reading in Data

---

.sav

```
library(haven)  
data <- read_sav((here("data", file-name.sav)))
```

# 1. Reading in Data

---

Google Sheets

```
library(googlesheets4)
```

Web

```
data <- read_csv("https://github.com/data-edu/dataedu/raw/master,
```

# 1. Reading in Data

---

Pitfalls

Can't find data file

- Check file location
- Check working directory
- Make sure those line up

# 1. Reading in Data

---

Other considerations to keep in mind

–Missing column name headers, csv files in particular –  
renaming column headers with `dplyr::rename()`

```
data <- readr::read_csv(here("data", "answer_export.csv"))
```

```
##  
## — Column specification —  
## cols(  
##   id = col_double(),  
##   content = col_character(),  
##   session_id = col_character(),  
##   question_id = col_double()  
## )
```

```
colnames(data)
```

```
## [1] "id"          "content"     "session_id"  "question_id"
```

```
data <- data %>%  
  rename(ID_num = id,
```

## 2. Tidying Data

---

Tidying data AKA wrangling data, cleaning data, etc.

- Data cleaning is 80% of the job
- Don't underestimate it, know what to expect

## 2. Tidying Data

---

"Happy families are all alike;  
every unhappy family is unhappy in its own way."  
-- Leo Tolstoy

"Tidy datasets are all alike,  
but every messy dataset is messy in its own way."  
-- Hadley Wickham

What does 'tidy' data mean

1. Every Column is a variable
2. Every Row is an observation
3. Every cell is a single value

(Don't worry too much about the details of what these mean quite yet)



## 2. Tidying data

---

Examining data is the first step

- `View()`
- `glimpse()`
- `colnames()` or `names()`
- `table()` or `tabyl()`

## 2. Tidying data

---

Deeper aspects of cleaning/tidying data

- Columns should be the right type (e.g. integer, double, factor, character, additional character types like date/time)
- Values in column should only be valid values for whatever that variable is supposed to be (sometimes your data sources will pre-validate this, sometimes not)
- Missing data indicators for missing values (NA type)

# 3. Our Tidy Data Tools

---

What are our tools for tidy data?

- Tidyverse
- More specifically tidyr and dplyr are the wrangling tools
- tidyr narrower focus on reshaping data
- dplyr does more, stuff like filtering arranging selecting joining
- dplyr also has powerful tools for grouping data and grouped operations
- stringr: not talked about as much but useful operations for dealing with character strings

# 3. Our Tidy Data Tools

---

How dplyr functions work:

- First argument is always data frame (important for %>%)
- They don't modify the existing data frame, need to save the results

```
data <- data %>%  
  rename(ID_num = id,  
         Response = content,  
         Session_num = session_id,  
         Question_num = question_id)
```

# 3. Our Tidy Data Tools

---

How dplyr functions work:

Very amenable to chained operations with `%>%`

```
data <- readr::read_csv(here("data", "answer_export.csv"))
```

```
##  
## — Column specification —————  
## cols(  
##   id = col_double(),  
##   content = col_character(),  
##   session_id = col_character(),  
##   question_id = col_double()  
## )
```

```
data_subset <- data %>%  
  rename(ID_num = id,  
         Response = content,  
         Session_num = session_id,  
         Question_num = question_id) %>%  
  select(Response, Question_num)
```

# 3. Our Tidy Data Tools

---

How dplyr functions work:

Grouped operations

```
data <- readr::read_csv(here("data", "answer_export.csv"))
```

```
##  
## — Column specification —————  
## cols(  
##   id = col_double(),  
##   content = col_character(),  
##   session_id = col_character(),  
##   question_id = col_double()  
## )
```

```
data_grp <- data_subset %>%  
  group_by(Question_num)
```

# Course Logistics

---

- Exam 1: postponed 1 week, will cover material through week 3
- Homework 3: Will be available by the end of the day tomorrow
- More demo videos on specific things, please share input on what you want to see on slack

# Wrapping up

---

In your base group's Slack channel:

- What is one thing you took away from today?
- What is something you want to learn more about?