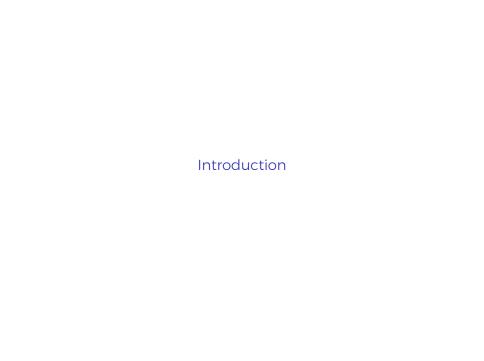
# Managing and Manipulating Data Using R

Lecture 8, Acquiring data in R

Ozan Jaquette



Libraries we will use today [install if you don't have them]

```
library(dplyr)
library(readr)
library(haven)
library(readxl)
library(labelled)
```

# Data we will use today

IPEDS or HSLS Federal Student Aid Chetty

Work on this



## Common data formats

- o Comma-separated values (.csv)
- Excel (.xls or .xlsx)
- Text-formated data (.txt)
- Tab-separated values (.tsv)
- o R (.Rdata or .rds)
- o Stata (.dta)
- o SPSS (.sav)
- o SAS (.sas)



#### readr

The readr package is part of tidyverse, which is designed to read in flat data files in R and transform them into data frames.

- We could load **library(tidyverse)** if we wanted to load all packages in tidyverse (e.g. ggplot2, dplyr, tidyr, stringr, readr, etc...)

- For the purpose of this lecture, we will just need to load library(readr)

readr

No matter the flat file format you are working with, there are two important steps for reading in data with  ${\tt readr}$ :

- 1. a function to parse the file (read\_csv)
- 2. column specification

Talk about readr (tidyverse), Notes from the reading, Notes from other sources

## readr functions

## readr's (tidyverse) functions

Format	Function			
Comma-separated values (csv)				
Semicolon separated	files			
Tab-separated values	(tsv)			
Any delimiter				
Fixed width files				
Text-formated data (t	xt)			
Web log files				

What type of data could be read in

## readr column specification

readr is pretty good at guessing each column's data type (e.g. character, double, etc.), however it is good practice to manually specify the data type for each column.

```
mtcars <- read_csv(readr_example("mtcars.csv"))</pre>
#> Parsed with column specification:
#> cols(
#> mpg = col double(),
#> cyl = col_integer(),
    disp = col double(),
#>
#> hp = col_integer(),
#> drat = col double(),
    wt = col double(),
#>
    gsec = col_double(),
#>
#> vs = col integer(),
#>
    am = col_integer(),
#>
    gear = col integer(),
     carb = col integer()
#>
#> )
```

## readr column specification

The output of the previous example shows us the column specification of the variables in the data. This is important because we could manually change column specification if we do not like readr's guess.

```
mtcars <- read_csv(readr_example("mtcars.csv"), col types =</pre>
  cols(
    mpg = col_double(),
    cyl = col_integer(),
   disp = col double(),
    hp = col_integer(),
    drat = col double(),
   vs = col_integer(),
    wt = col_double(),
    qsec = col_double(),
    am = col_integer(),
    gear = col_integer(),
    carb = col_integer()
```

## readr features

skip: comment: col\_names:

readr automatically treats the first line of data as column names.

There are instances where you may want to tell R from what line to begin reading in data.

Notice the example below. The first two lines are comments about the data. We would need to use  $\mathbf{skip} = \mathbf{n}$  to skip n lines.

We could also tell R to drop lines we specify as comments. With comment = n

#### readr column names

We could tell R there are no column names with **col\_names = FALSE** or we could manually give R column names with **col\_names = c("", "", "")** 

#### readr Student exercise

- Get in your homework groups
- Create a 3x3 tibble like the examples above (e.g. read\_csv("a,b,c....")), treating the first line as column names
- Now on the first line add a sentence
- This time add a special character (\*, #,!) at the beginning of the sentence and indicate it is a comment
- Delete the sentence and column names (should have a 2x2 tibble) and manually tell R column names

# NOT SURE IF TO MAKE THIS A DEMONSTRATION WHERE STUDENTS FOLLOW ALONG OR ANOTHER STUDENT EXERCISE Tying it all together

Use read\_csv() function from readr to import csv dataset into R without column specification. Follow along on your computers.

```
#ipeds_ay <- read_csv(file="../ic2017ay_small.csv")</pre>
```

Notice tuition and fee columns are read in as character type.

Use read\_csv() function from readr to import csv dataset into R with column specification [Would it be better to change to integer or double?]

```
# {ipeds_ay <- read_csv("../ic2017ay_small.csv", col_types =
  cols(
    unitid = col_integer(),
    tuition1 = col integer(),
    fee1 = col integer(),
    hrchg1 = col_integer(),
    tuition2 = col integer(),
    fee2 = col_integer(),
    hrchg2 = col_integer(),
    tuition3 = col integer(),
    fee3 = col_integer(),
    hrchg3 = col integer()
# )
#) # }
```

## readr Running into errors

- 1. Make sure you have downloaded and saved flat file
- Make sure to know the file path of where data is downloaded or saved (~/Desktop/educ263/data)
- Make sure you set your working setwd() directory in R. To check your current working directory type getwd() in console.



#### haven

## Recap from lecture 5

haven is part of **tidyverse**, which enables users to import and export data from the following statistical packages:

- SAS
- SPSS
- Stata

Similar to readr, we could load the entire **library(tidyverse)** package to get haven. For the purpose of this lecture, we will just need to load **library(haven)**.

## haven functions

haven's (tidyverse) functions

Format	Function
SPSS	read_sav
SAS	read_sas
Stata	read_dta

# haven read and write Stata arguments

```
read_dta(file, encoding = NULL)
write data(data, path, version = 14)
```

#### **Arguments**

- file: file path to data
- encoding: files prior to Stata 14 did not declare text encoding, files after Stata 14 do not need to declare encoding value
- data: data frame to save (write)
- path: file path to where data will be saved
- version: file version

Link

## haven read and write Stata example

Use read\_dta() function from haven to import Stata dataset into R

```
hsls <- read_dta("~/Desktop/lecture8/hsls_sch_small.dta", encoding=NULL)

# View data
head(hsls)
glimpse(hsls)
```

Use write\_dta function from haven to write State data

```
write_dta(hsls, path = "~/Desktop/lecture8/hsls_sch_small.dta")
```

Student exercise Running into problems with variable/value labels?



#### readxl

The  ${\tt readxl}$  package is part of tidyverse, which is designed to easily read data from Excel and into R.

- We could load **library(tidyverse)** if we wanted to load all packages in tidyverse. For the purpose of this lecture, we just need to load **library(readxl)**.

#### readxl

readx1 supports both .xls and .xlsx formats and is designed to work with tabular data. It does not require dependencies- making installing and operating fairly simple.

readx1 has several example files where we could use as practice. The files include:

For now, lets use "datasets.xlsx"

```
excel_example <- readxl_example("datasets.xlsx")</pre>
```

## readxl features

- sheet: read excel(excel file, sheet = "sheet name")
- o n\_max: read\_excel(excel file, n max = n)
- o range: read\_excel(excel file, range = "A:D")
- o cell\_rows: read\_excel(excel file, range = cell\_rows(1:n))
- o cell\_cols: read\_excel(excel file, range = cell\_cols("A:D"))
- o na: read\_excel(excel file, na = "n")

#### readxl sheet

```
#To view sheets in excel file
excel_sheets(excel_example)
#> [1] "iris" "mtcars" "chickwts" "quakes"
xl_example <- read_excel(excel_example, sheet = "quakes")</pre>
head(xl example)
#> # A tibble: 6 x 5
#> lat long depth mag stations
#> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 -20.42 181.62 562 4.8
                               41
#> 2 -20.62 181.03 650 4.2 15
#> 3 -26.00 184.10 42 5.4 43
#> 4 -17.97 181.66 626 4.1 19
#> 5 -20.42 181.96
                  649 4.0 11
#> 6 -19.68 184.31 195 4.0 12
```

## readxl n\_max

```
read_excel(excel_example, sheet = "quakes", n_max = 3)
#> # A tibble: 3 x 5
#> lat long depth mag stations
#> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
#> 1 -20.42 181.62 562 4.8 41
#> 2 -20.62 181.03 650 4.2 15
#> 3 -26.00 184.10 42 5.4 43
```

## readxl range

```
read_excel(excel example, sheet = "quakes", range = "C1:E4")
#> # A tibble: 3 x 3
#> depth mag stations
#> <dbl> <dbl> <dbl>
#> 1 562 4.8 41
#> 2 650 4.2 15
#> 3 42 5.4 43
read_excel(excel_example, sheet = "quakes", range = cell_rows(1:3))
#> # A tibble: 2 x 5
#> lat long depth mag stations
#> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 -20.42 181.62 562 4.8 41
#> 2 -20.62 181.03 650 4.2 15
head(read_excel(excel_example, sheet = "quakes", range = cell_cols("A:C")))
#> # A tibble: 6 x 3
#> lat long depth
#> <db1> <db1> <db1>
#> 1 -20.42 181.62 562
#> 2 -20.62 181.03 650
#> 3 -26.00 184.10 42
#> 4 -17.97 181.66 626
#> 5 -20.42 181.96 649
#> 6 -19.68 184.31 195
# using head() to only view first 6 rows
```

#### readxl na

```
read_excel(excel_example, sheet = "quakes", na = "-20.42")
#> # A tibble: 1,000 x 5
     lat long depth mag stations
#>
#> <dbl> <dbl> <dbl> <dbl> <dbl>
       NA 181.62 562 4.8 41
#> 1
#> 2 -20.62 181.03 650 4.2 15
#> 3 -26.00 184.10 42 5.4 43
#> 4 -17.97 181.66 626 4.1 19
#> 5 NA 181.96 649 4.0
                         11
#> 6 -19.68 184.31 195 4.0
                         12
#> 7 -11.70 166.10 82 4.8 43
#> 8 -28.11 181.93 194 4.4 15
#> 9 -28.74 181.74 211 4.7 35
#> 10 -17.47 179.59 622 4.3 19
#> # ... with 990 more rows
```

## readxl Student exercise

Save data - Download and save Federal Student Financial Aid Data

- Read in data using readx1 function
- Read in first four rows (n $\_$ max) Read in column Names to column State **hint** cell $\_$ cols Set value "A" to missing (na) **note** : you need to investigate in detail before setting anything to missing

# readxl Running into problems

- 1. Make sure you have downloaded and saved excel file
- Make sure to know the file path of where data is downloaded or saved (~/Desktop/educ263/data)
- Make sure you set your working setwd() directory in R. To check your current working directory type getwd() in console.
- 4. Make sure to choose the correct sheet (if applicable)
- 5. Pay attention to column names when setting range



## Downloading data from web

Could save time and reduce the steps of downloading, saving, and reading in data, by reading in data directly from the web.

- note that not all packages will work downloading data from web (read excel)

For example, rather than downloading ipeds data and saving it in a folder, we could download the data directly from the web.

# Downloading data from web example using Raj Chetty data

- Follow this link and under the "Mobility Report Cards..." tab select "click to view data".
- 2. Choose "Online Data Table 1"
- 3. Right click and copy link address for "Excel" (Note: it is actually a csv file)

#### Mobility Report Cards: The Role of Colleges in Intergenerational Mobility

Chetty, Friedman, Saez, Turner, and Yagan (2017)

Mobility Statistics and Student Outcomes by College and Birth Cohort

Click to view data -

Data Description		Download		
Online Data Table 1 Preferred Estimates of Access and Mobility Rates by College	Stata	Excel	Readme	
Online Data Table 2 Baseline Cross-Sectional Estimates by College	Stata	Excel	Readme	
Online Data Table 3 Baseline Longitudinal Estimates by College and Child's Cohort	Stata	Excel	Readme	

## Downloading data from web

#### [FIX OUTPUT (CUTTING OFF)]

```
#Paste url to excel "csv" file
data url <- "http://www.equality-of-opportunity.org/data/college/mrc table1.csv"
#Download data and read in using read_csv (readr)
mrc <- read csv(data url)
#View first 4 rows and 4 columns
mrc[1:4, 1:4]
#> # A tibble: 4 x 4
#> super opeid
                                                      name czname state
#>
          <int.>
                                                     <chr> <chr> <chr> <chr>
#> 1
          2665 Vaughn College Of Aeronautics And Technology New York NY
#> 2
          7273
                             CUNY Bernard M. Baruch College New York NY
                            City College Of New York - CUNY New York NY
#> 3
           2688
          7022
                                       CUNY Lehman College New York
#> 4
                                                                      NY
```

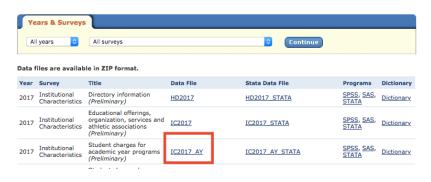
## Downloading data from web

## Alternative approach

```
#Download data and read in link directly using read csv (readr)
mrc <- read_csv("http://www.equality-of-opportunity.org/data/college/mrc table1.
#> Parsed with column specification:
#> cols(
#>
    super opeid = col integer(),
#>
    name = col_character(),
    czname = col character(),
    state = col character(),
#>
#>
    par_median = col_integer(),
    k median = col integer(),
#>
#>
     par_q1 = col_double(),
#>
    par_top1pc = col_double(),
    kg5 cond parg1 = col double(),
#>
     ktop1pc_cond_parq1 = col_double(),
#>
     mr_kq5_pq1 = col_double(),
#>
    mr_ktop1_pq1 = col_double(),
#>
#>
     trend_parq1 = col_double(),
     trend bottom40 = col double(),
#>
#>
     count = col double()
#> )
#View first 4 rows and 4 columns
mrc[1:4, 1:4]
#> # A tibble: 4 x 4
#> super opeid
                                                                 czname state
                                                          name
```

## Problems downloading data (zip files) using IPEDS

- 1. Follow this link and under the "Survey Data" tab select "Complete data files".
- 2. Choose "All years" and "All surveys" and click continue
- 3. Right click and copy link address for "IC2017\_AY"



# Downloading data (zip files) using IPEDS

```
# Paste url and read in using read_csv
# What happens when you try reading in this zip file?
#ic2017_ay <- read_csv("https://nces.ed.gov/ipeds/datacenter/data/IC2017_AY.zip"</pre>
```

# Need to download file first to unzip