## Data Imports

Fall 2022

October 03 2022

## Working Directories

The working directory is where R looks for files and saves files by default.

```
getwd() # see working directory
setwd() # change your working directory
```

#### To set working directory to your STAT 220 course folder

```
setwd("path/to/stat220-folder/") # set
getwd() # check
```

#### **Useful Terminal Commands:**

```
$ cd  # change directory
$ ls  # unix command to list files
$ pwd  # present working directory
```

## Web Imports

#### To your working environment:

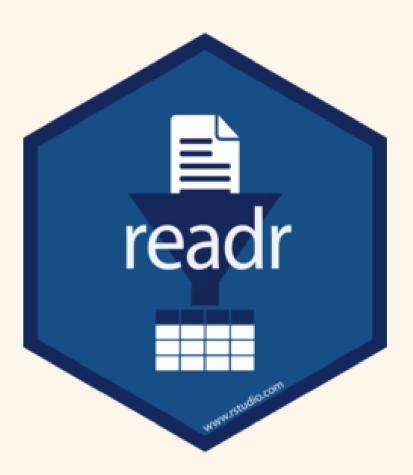
```
url <- "https://raw.githubusercontent.com/deepbas/statdatasets/main/murders.csv"
dat <- read_csv(url)</pre>
```

#### To download file to working folder:

```
download.file(url, "murders.csv")
```

#### readr

- readr is a part of tidyverse library
- Includes functions for reading data stored in text file spreadsheets into R.
- Functions in the package include read\_csv(), read\_tsv(), read\_delim() and more.
- These differ by the delimiter they use to split columns.



## readr functions

function	reads
read_csv()	Comma separated values
read_csv2()	Semi-colon separated values
read_delim()	General delimited files
read_fwf()	Fixed width files
read_log()	Apache log files
read_table()	Space separated
read_tsv()	Tab delimited values

## Basic syntax

All readr functions share a common syntax

```
df <- read_csv(file = "path/to/file.csv", ...)</pre>
```

## Base R Imports

- R-base import functions
  - o read.csv()
  - o read.table()
  - o read.delim()
- Generate data frames rather than tibbles
- Character variables are converted to factors
  - Can be avoided by setting the argument stringsAsFactors=FALSE

## Advantages of readr

#### readr functions are:

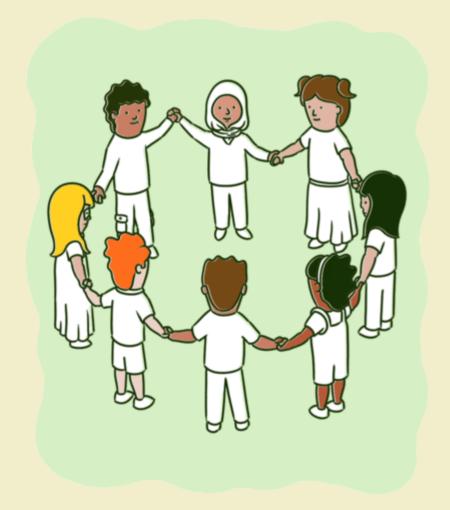
- ~ 10 times faster
- Return tibbles
- Have more intuitive defaults.
- No row names, no strings as factors.

### Data frames and tibbles Conversion



- as\_tibble() convert a data frame to a tibble
- as.data.frame() convert a tibble to a data frame

# Group Activity 1



- Let's go over to maize server/ local Rstudio and our class moodle
- Get the class activity 10.Rmd file
- Skim through problem 1

## Did it work as expected?

```
Rows: 549
Columns: 16
$ series
                     $ episode
                     <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, ...
$ baker
                     <chr> "Annetha", "David", "Edd", "Jasminder", "Jonatha...
                     <chr> "2nd", "3rd", "1st", "N/A", "9th", "N/A", "8th",...
$ technical
$ result
                     <chr> "IN", "IN", "IN", "IN", "IN", "IN", "IN", "IN", "IN", ...
$ uk airdate
                     <chr> "17 August 2010", "17 August 2010", "17 August 2...
$ us season
                     $ us airdate
                     $ showstopper_chocolate <chr>> "chocolate", "chocolate", "no chocolate", "no ch...
$ showstopper dessert
                     <chr> "other", "other", "other", "other", "ca...
$ showstopper fruit
                     <chr> "no fruit", "no fruit", "no fruit", "no fruit", ...
$ showstopper nut
                     <chr> "no nut", "no nut", "no nut", "no nut", "almond"...
                     <chr> "no chocolate", "chocolate", "no chocolate", "no...
$ signature_chocolate
$ signature_dessert
                     <chr> "cake", "cake", "cake", "cake", "cake", "cake", ...
$ signature_fruit
                     <chr> "no fruit", "fruit", "fruit", "fruit", "fruit", ...
$ signature nut
                     <chr> "no nut", "no nut", "no nut", "no nut", "no nut"...
```

We want technical to be numerical and uk\_airdate to be date

## The col\_types argument

By default, looks at first 1000 rows to guess variable data types (guess\_max)

```
desserts <- read_csv(
  "https://raw.githubusercontent.com/deepbas/statdatasets/main/desserts.csv",
  col_types = list(
    technical = col_number(),
    uk_airdate = col_date()
  )
)</pre>
```

## Looking for problems

#### List of potential problems parsing the file

```
problems(desserts)
# A tibble: 556 × 5
           col expected
                                actual
                                                file
     row
   <int> <int> <chr>
                                <chr>
                                                <chr>
             6 date in IS08601 17 August 2010
             6 date in ISO8601 17 August 2010
 3
             6 date in ISO8601 17 August 2010
             4 a number
                                N/A
                                                 11 11
             6 date in ISO8601 17 August 2010
 6
             6 date in ISO8601 17 August 2010
                                                 11.11
             4 a number
                                N/A
             6 date in ISO8601 17 August 2010
             6 date in ISO8601 17 August 2010
10
             4 a number
                                N/A
                                                 11.11
# ... with 546 more rows
```

## Date formatting

```
# A tibble: 556 × 5
          col expected actual
                                               file
    row
  <int> <int> <chr>
                              <chr>
                                                <chr>
            6 date in ISO8601 17 August 2010 ""
     3 6 date in IS
4 6 date in IS
5 4 a number
2
3
4
            6 date in ISO8601 17 August 2010
            6 date in ISO8601 17 August 2010
                                                11 11
                               N/A
      5 6 date in ISO8601 17 August 2010
# ... with 551 more rows
```

ISO8601 format: 2021-10-04

What we have: 17 August 2010

## Adding format instructions

```
desserts <- read_csv(
  "https://raw.githubusercontent.com/deepbas/statdatasets/main/desserts.csv",
    col_types = list(
    technical = col_number(),
    uk_airdate = col_date(format = "%d %B %Y")
  )
)</pre>
```

- Year: "%Y" (4 digits). "%y" (2 digits)
- Month: "%m" (2 digits), "%b" (abbreviated name in current locale), "%B" (full name in current locale).
- Day: "%d" (2 digits), "%e" (optional leading space)

## Looking for more problems

#### List of potential problems parsing the file

```
problems(desserts)
# A tibble: 7 \times 5
        col expected actual file
    row
  <int> <int> <chr> <chr>
            4 a number N/A
                               11.11
      5
2
            4 a number N/A
                               11.11
      9 4 a number N/A
                               11.11
     11 4 a number N/A
                               11.11
     4 a number N/A
                              11.11
     36 4 a number N/A
                               11 11
6
            4 a number N/A
                               11.11
```

## Addressing missing values

By default na = c("", "NA") are the recognized missing values

```
desserts <- read_csv(
   "https://raw.githubusercontent.com/deepbas/statdatasets/main/desserts.csv",
   col_types = list(
     technical = col_number(),
     uk_airdate = col_date(format = "%d %B %Y")
   ),
   na = c("", "NA", "N/A")
)</pre>
```

#### No more problems

```
problems(desserts)
# A tibble: 0 × 5
# ... with 5 variables: row <int>, col <int>, expected <chr>, actual <chr>,
# file <chr>
```

#### The Dataset

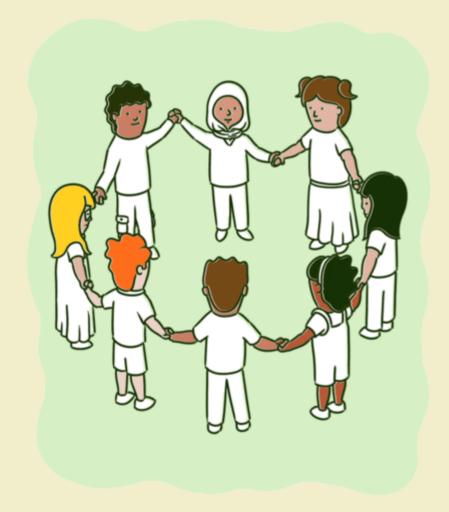
```
# A tibble: 549 × 16
   series episode baker
                          techn...¹ result uk airdate us se...² us airdate shows...³
    <fdb>
             <dbl> <chr>
                                <dbl> <chr>
                                               <date>
                                                             <dbl> <date>
                                                                                <chr>
                 1 Annetha
                                     2 IN
                                               2010-08-17
                                                                NA NA
                                                                                chocol...
                 1 David
                                     3 TN
                                               2010-08-17
                                                                NA NA
                                                                                chocol...
 3
                 1 Edd
                                     1 IN
                                               2010-08-17
                                                                NA NA
                                                                                no cho...
 4
                 1 Jasminder
                                    NA IN
                                               2010-08-17
                                                                NA NA
                                                                                no cho...
                 1 Jonathan
                                     9 IN
                                               2010-08-17
                                                                NA NA
                                                                                no cho...
 6
                 1 Louise
                                    NA IN
                                               2010-08-17
                                                                                chocol...
                                                                NA NA
                 1 Miranda
                                     8 IN
                                               2010-08-17
                                                                NA NA
                                                                                chocol...
 8
                 1 Ruth
                                    NA IN
                                               2010-08-17
                                                                NA NA
                                                                                chocol...
 9
                                    10 OUT
                                                                                chocol...
                 1 Lea
                                               2010-08-17
                                                                NA NA
10
                                    NA OUT
                                               2010-08-17
                                                                                chocol...
                 1 Mark
                                                                NA NA
# ... with 539 more rows, 7 more variables: showstopper_dessert <chr>,
#
    showstopper_fruit <chr>, showstopper_nut <chr>, signature_chocolate <chr>,
    signature_dessert <chr>, signature_fruit <chr>, signature_nut <chr>, and
    abbreviated variable names <sup>1</sup>technical, <sup>2</sup>us_season, <sup>3</sup>showstopper_chocolate
#
```

## Column casting functions

Туре	<pre>dplyr::glimpse()</pre>	readr::col_*()
logical	<lgl></lgl>	col_logical
numeric	<int> or <dbl></dbl></int>	col_number
character	<chr></chr>	col_character
factor	<fct></fct>	col_factor
date	<date></date>	col_date

## ?read\_csv

# Group Activity 2



- Work on problem 2 to fix some messy data
- Ask me questions

### Another example

```
energy %>% slice_max(n=2, order_by = Timestamp)
# A tibble: 2 × 90
 Timestamp year month weekOfYear dayOfMonth dayWeek timeHour timeM...<sup>1</sup>
                   <dbl> <dbl> <dbl> <dbl> <fct> <dbl>
 <dttm>
                                                                      <dbl>
1 2016-08-31 23:45:00 2016 8
                                       35
                                                 31 Wed
                                                                 23
                                                                         45
                                                 31 Wed
2 2016-08-31 23:30:00 2016 8
                                      35
                                                                 23
                                                                         30
# ... with 82 more variables: `100 Nevada Street` <dbl>, `104 Maple St.` <dbl>,
   `106_Winona_St.` <dbl>, Allen_House <dbl>,
   `Alumni_Guest_House/Johnson_House` <dbl>, Arboretum_Office <dbl>,
   Art_Studios <dbl>, Benton_House <dbl>, Berg_House <dbl>, Bird_House <dbl>,
#
   Boliou_Memorial_Art_Bldg. <dbl>, Burton_Hall <dbl>,
#
   `Cassat_Hall_/_James_Hall` <dbl>,
#
   `Center_for_Mathematics_&_Computing` <dbl>, Chaney_House <dbl>, ...
```

## Wide to Long

```
energy_narrow <- energy %>%
  pivot_longer(names_to = "building", values_to = "energyKWH", cols = `100_Nevada_Street`:Wi
```

```
energy narrow
# A tibble: 2,880,578 × 10
   Timestamp
             year month weekOfYear dayOfMonth dayWeek timeH...¹ timeM...²
                     <dbl> <dbl>
                                        <dbl>
                                              <dbl> <fct> <dbl> <dbl>
   <dttm>
 1 2015-09-01 00:00:00 2015
                                           35
                                                        1 Tues
 2 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
 3 2015-09-01 00:00:00
                        2015
                                           35
                                                        1 Tues
 4 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
 5 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
 6 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
 7 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
 8 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
 9 2015-09-01 00:00:00
                        2015
                                           35
                                                       1 Tues
10 2015-09-01 00:00:00 2015
                                           35
                                                        1 Tues
# ... with 2,880,568 more rows, 2 more variables: building <chr>,
    energyKWH <dbl>, and abbreviated variable names <sup>1</sup>timeHour, <sup>2</sup>timeMinute
```

## More dates and times manipulation

```
energy$Timestamp[1]
[1] "2015-09-01 UTC"
## [1] "2015-09-01 UTC"
as.numeric(energy$Timestamp[1])
[1] 1441065600
## [1] 1441065600

# 5th timestamp
( stamp5 <- energy$Timestamp[5] )
[1] "2015-09-01 01:00:00 UTC"</pre>
```

```
mdy("1/4/2021")
[1] "2021-01-04"

dmy_hms("01/04/2020-01-30-23")
[1] "2020-04-01 01:30:23 UTC"
```

#### Another example: duration using lubriate

```
top_dest <- flights %>%
 count(dest) %>%
  slice max(n, n = 10)
flights %>%
  semi_join(top_dest) %>%
  mutate(sch_datetime = make_datetime(
    year = year, month = month,
    day = day, hour = hour,
   min = minute)
    ) %>%
  select(dest, sch_datetime) %>%
  group by(dest) %>%
  arrange(sch_datetime) %>%
  mutate(
    diff1 = (lag(sch_datetime) %--%
                    sch datetime)/minutes(1),
    diff2 = interval(lag(sch_datetime),
                    sch datetime)/minutes(1)) %>%
  summarize(medianMins1 = median(diff1,
                                 na.rm=TRUE),
            medianMins2 = median(diff2,
                                 na.rm=TRUE))
```

```
# A tibble: 10 \times 3
   dest medianMins1 medianMins2
                <dbl>
                              <dbl>
   <chr>
 1 ATL
                    15
                                 15
 2 BOS
                    17
                                 17
 3 CLT
                    18
                                 18
 4 DCA
                   34
                                 34
 5 FLL
                   24
                                 24
 6 LAX
                    19
                                 19
 7 MCO
                    20
                                 20
 8 MIA
                   25
                                 25
 9 ORD
                   15
                                 15
10 SFO
                    20
                                 20
```