tidyverse

Matt McDonald

inspired by

https://github.com/michaellevy/tidyverse_talk/blob/master/tidyverse.md

What is the tidyverse?

- R packages for data science
- The tidyverse is an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

Tidy data

Put data in data frames

- Each type of observation gets a data frame
- Each variable gets a column
- Each observation gets a row

Tidy APIs

Functions should be consistent and easily (human) readable

- Take one step at a time
- Connect simple steps with the pipe
- Referential transparency

Okay but really, what is it?

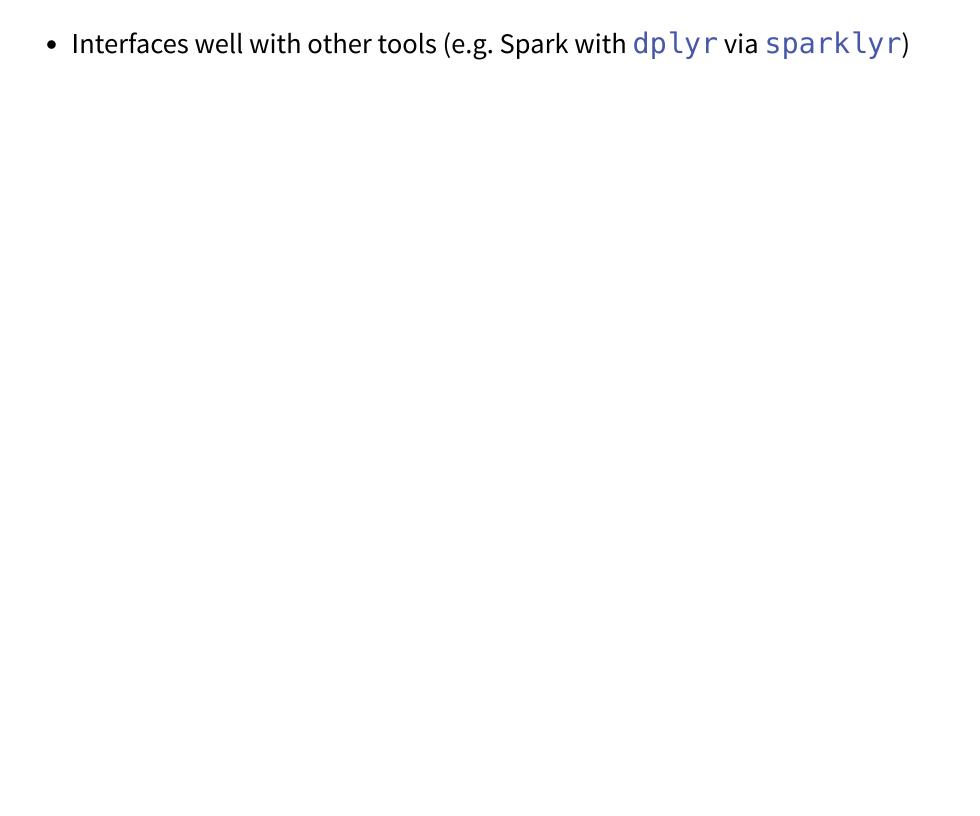
- Suite of ~20 packages that provide consistent, user-friendly, smart-default tools to do most of what most people do in R.
- https://www.tidyverse.org/packages/
- install.packages (tidyverse) installs all of the above packages.
- library(tidyverse) attaches only the core packages.

Why tidyverse?

- Consistency
 - e.g. All stringr functions take string first
 - e.g. Many functions take data.frame first -> piping
 - Faster to write
 - Easier to read
 - Tidy data: Imposes good practices
 - Type specificity
- Implements simple solutions to common problems (e.g. purrr::transpose)
- Smarter defaults

```
e.g. utils::write.csv(row.names = FALSE) =
readr::write_csv()
```

Runs fast (thanks to Rcpp)



Data Types in R

- **Numeric**: This is the default data type for numbers in R. It includes real numbers (floating-point values) and integers. For example, 42, 3.14.
- Integer: Specifically for integer values. While numeric data can include integers, if you specifically want to declare an integer, you append L to the number, like 42L.
- Logical: This type represents boolean values and can either be TRUE or FALSE.
- **Character**: This type represents strings. Text and characters are enclosed in quotes. For example, "Hello, World!".
- **Factor**: A data type used for categorical data. Factors can be ordered or unordered and are very useful in statistical modeling and graphics. They are stored as integers but each integer value corresponds to a label.
- Other: Complex, Raw and Date.

Data Structures in R

- **Vectors**: An ordered collection of elements of the same basic data type.
- **Matrices**: Two-dimensional, rectangular layouts of elements of the same basic data type.
- Arrays: Similar to matrices but can have more than two dimensions.
- **Data frames**: A table or a two-dimensional array-like structure where each column can contain different types of data (numeric, character, factor, etc.). It's one of the most important data types in R for data analysis.
- **Lists**: An ordered collection of objects (components). A list in R can contain objects of different types including numbers, strings, vectors, and even other lists.

tibble

Tibbles are a modern re-imagining of data frames.

Tibbles print politely.

```
tdf = tibble(x = 1:1e4, y = rnorm(1e4))
tdf
```

 Can customize print methods with print(tdf,

```
n = rows, width =
cols)
```

• Set default with

```
options(tibble.print_max
= rows, tibble.width
= cols)
```

Tibble defaults

Tibbles have some convenient and consistent defaults that are different from base R data.frames.

type consistency

```
dfs = list(
    df = data.frame(abc = letters[1:3], xyz = letters[24:2
    tbl = tibble(abc = letters[1:3], xyz = letters[24:26])
)

sapply(dfs, function(d) class(d[, "abc"]))

$df
[1] "character"

$tbl
[1] "tbl_df" "tbl" "data.frame"
```

Note that tidyverse import functions (e.g. readr::read_csv) default to tibbles and that this can break existing code.

List-columns!

```
a <- tibble(ints = 1:5,</pre>
            powers = lapply(1:5, function(x) x^{(1:x)})
    a[[5,2]]
[[1]]
[1]
       5 25 125 625 3125
    a
# A tibble: 5 \times 2
   ints powers
  <int> <list>
      1 <dbl [1]>
2 2 <dbl [2]>
3 3 <dbl [3]>
4 4 <dbl [4]>
      5 <dbl [5]>
```

The pipe %>%

Sends the output of the LHS function to the first argument of the RHS function.

```
1:8 %>%
    sum() %>%
    sqrt()

[1] 6

sqrt(sum(1:8))

[1] 6
```

dplyr

Common data(frame) manipulation tasks. Four core "verbs":

- filter
- select
- arrange
- group_by + summarize

We will cover in more detail in further classes

joins

dplyr also does multi-table joins and can connect to various types of databases.

```
t1 = tibble(alpha = c(letters[2:7],letters[11]), num = c(2:7,11))
t2 = tibble(alpha = letters[4:10], num = 4:10)
t3 <- full_join(t1, t2, by = "alpha", suffix = c("_t1", "_t2"))
arrange(t3, alpha)</pre>
```

```
# A tibble: 10 \times 3
   alpha num_t1 num_t2
   <chr> <dbl> <int>
 1 b
                     NA
 2 c
                     NA
 3 d
 4 e
                      5
 5 f
 6 g
 7 h
             NA
8 i
                      9
             NA
9 j
                     10
             NA
10 k
             11
                     NA
```

Other Useful Packages

- ggplot2 data visualization
- tidyr data reshaping
- *stringr* manipulating strings
- *purrr* programming + working with lists

What does "un-tidy" data mean?

```
who <- read csv(here::here('Lecture1', 'who.csv'))</pre>
    who
# A tibble: 9,137 \times 46
   country iso2 iso3
                          year new sp m014 new sp m1524 new sp m2534 new sp m3544
            <chr> <chr> <dbl>
                                      <dbl>
                                                    <dbl>
                                                                  <dbl>
   <chr>
                                                                                 <dbl>
 1 Afghani... AF
                   AFG
                           1980
                                         NA
                                                        NA
                                                                     NA
                                                                                   NA
 2 Afghani... AF
                   AFG
                           1981
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
                          1982
 3 Afghani... AF
                   AFG
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
                          1983
 4 Afghani... AF
                   AFG
                                                                                   NA
                                         NA
                                                       NA
                                                                     NA
 5 Afghani... AF
                          1984
                   AFG
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
 6 Afghani... AF
                   AFG
                          1985
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
 7 Afghani... AF
                   AFG
                          1986
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
8 Afghani... AF
                   AFG
                          1987
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
 9 Afghani... AF
                   AFG
                          1988
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
10 Afghani... AF
                   AFG
                          1989
                                         NA
                                                       NA
                                                                     NA
                                                                                   NA
# i 9,127 more rows
# i 38 more variables: new sp m4554 <dbl>, new sp m5564 <dbl>,
    new sp m65 <dbl>, new sp f014 <dbl>, new sp f1524 <dbl>,
#
    new_sp_f2534 <dbl>, new_sp_f3544 <dbl>, new_sp_f4554 <dbl>,
#
    new sp f5564 <dbl>, new sp f65 <dbl>, new sn m014 <dbl>,
#
    new sn m1524 <dbl>, new sn m2534 <dbl>, new sn m3544 <dbl>,
```

new sn m4554 <dbl>, new sn m5564 <dbl>, new sn m65 <dbl>, ...

Using tidyverse to work with "un-tidy" data

```
who %>%
  select(-iso2, -iso3) %>%
  gather(group, cases, -country, -year ) %>%
  mutate(group = str_replace(group, "new_*", ""),
         method = str extract(group, "[a-z]+"),
         gender = str sub(str extract(group, " [a-z]"), 2, 2),
         age = str_extract(group, "[0-9]+"),
         age = ifelse(str length(age) > 2,
                      str c(str sub(age, 1, -3), str sub(age, -2, -1), sep = "-")
                      str c(age, "+"))) %>%
  group_by(year, gender, age, method) %>%
  summarize(total cases = sum(cases, na.rm = TRUE), .groups='drop') %>%
  ggplot(aes(x = year, y = total cases, linetype = gender)) +
  geom line() +
  facet grid(method ~ age,
             labeller = labeller(.rows = label both, .cols = label both)) +
  scale_y_log10() +
  theme light() +
  theme(axis.text.x = element text(angle = 45, vjust = 1, hjust = 1))
```

Using tidyverse to work with "un-tidy" data

