# Intro to R

Manipulating Data in R

#### Recap of Subsetting

- pull() to get values out of a data frame
- select() creates a smaller data frame with only certain columns
- you can select() based on patterns in the column names using tidyselect functions
- you can combine multiple tidyselect functions together like select(starts\_with("C"), ends\_with("state"))
- filter() can be used to filter out rows based on logical conditions
- avoid using quotes when referring to column names with filter()

### **Recap Subsetting Continued**

- $\cdot$  == is the same as equivalent to
- & means both conditions must be met to remain after filter()
- I means either conditions needs to be met to remain after filter()

#### Recap of Data Cleaning

- count() can help determine if we have NA values
- filter() automatically removes NA values can't confirm or deny if condition is met (need | is.na() to keep them)
- drop\_na() can help you remove NA values from a variable or an entire data frame
- NA values can change your calculation results
- think about what NA values represent

### Recap of Data Cleaning

 recode() can help with simple recoding values of a variable (needs to be inside the mutate function)

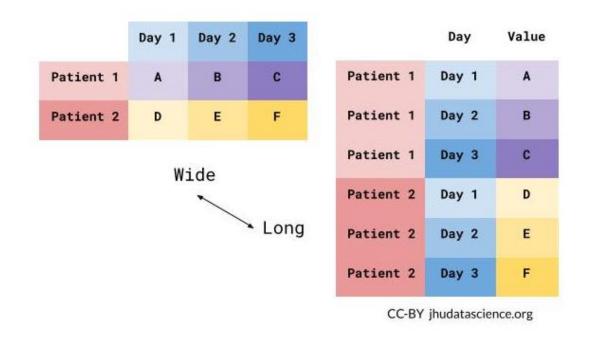
Cheatsheet

#### **Manipulating Data**

In this module, we will show you how to:

- 1. The two major forms of data (long and wide)
- 2. Reshape data from wide to long for data analysis and visualization

# Data is wide or long with respect to certain variables.



Data is stored *differently* in the tibble.

Wide: has many columns

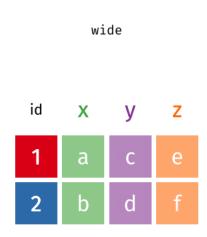
Long: column names become data

Wide: multiple columns per individual, values spread across multiple columns

Long: multiple rows per observation, a single column contains the values

```
# A tibble: 6 \times 3
  State
                          value
         name
  <chr> <chr>
                          <dbl>
                         0.516
1 Alabama June vacc rate
2 Alabama May_vacc_rate
                          0.514
3 Alabama April_vacc_rate 0.511
4 Alaska June vacc rate
                          0.627
5 Alaska May vacc rate
                          0.626
6 Alaska April vacc rate 0.623
```

https://github.com/gadenbuie/tidyexplain/blob/main/images/tidyr-pivoting.gif



#### Why do we need to switch between wide/long data?

#### Wide: Easier for humans to read

#### Long: Easier for R to make plots & do analysis

```
# A tibble: 6 \times 3
                          value
  State
         name
  <chr> <chr>
                          <dbl>
1 Alabama June vacc rate
                          0.516
2 Alabama May_vacc_rate
                          0.514
3 Alabama April_vacc_rate 0.511
4 Alaska June vacc rate
                          0.627
5 Alaska May vacc rate
                          0.626
6 Alaska April_vacc_rate 0.623
```

### Pivoting using tidyr package

tidyr allows you to "tidy" your data. We will be talking about:

pivot\_longer - make multiple columns into variables, (wide to long)

pivot\_longer...

pivot\_longer() - puts column data into rows (tidyr package)

First describe which columns we want to "pivot\_longer"

{long\_data} <- {wide\_data} %>% pivot\_longer(cols = {columns to pivot})

#### long\_data

wide\_data

pivot\_longer() - puts column data into rows (tidyr package)

- First describe which columns we want to "pivot\_longer"
- names\_to = gives a new name to the pivoted columns
- values\_to = gives a new name to the values that used to be in those columns

2 May\_vacc\_rate 0.514 3 April vacc\_rate 0.511

```
wide_data
# A tibble: 1 \times 3
  June_vacc_rate May_vacc_rate April_vacc_rate
          <dbl> <dbl>
                                        <fdb>>
1
          0.516 0.514
                                        0.511
long_data <- wide_data %>% pivot_longer(cols = everything(),
                                       names_to = "Month",
values_to = "Rate")
long data
# A tibble: 3 \times 2
 Month
                  Rate
 <chr>
        <dbl>
1 June_vacc_rate 0.516
```

Newly created column names are enclosed in quotation marks.

#### Data used: Charm City Circulator

https://hutchdatascience.org/SeattleStatSummer\_R/data/Charm\_City\_Circulator\_Rider

circ <- read\_csv("https://hutchdatascience.org/SeattleStatSummer\_R/data/Charm\_ head(circ, 5)

```
# A tibble: 5 \times 15
       date orangeBoardings orangeAlightings orangeAverage purpleBoardir
  day
  <chr> <chr>
                            <dbl>
                                             <db1>
                                                            <dbl>
1 Monday 01/1...
                              877
                                               1027
                                                             952
2 Tuesday 01/1...
                              777
                                               815
                                                             796
3 Wednesday 01/1...
                             1203
                                              1220
                                                            1212.
4 Thursday 01/1...
                             1194
                                              1233
                                                            1214.
5 Friday
            01/1...
                             1645
                                                            1644
                                               1643
# ... with 9 more variables: purpleAlightings <dbl>, purpleAverage <dbl>,
    greenBoardings <dbl>, greenAlightings <dbl>, greenAverage <dbl>,
#
    bannerBoardings <dbl>, bannerAlightings <dbl>, bannerAverage <dbl>,
    daily <dbl>
```

#### Mission: Taking the averages by line

Let's imagine we want to create a table of average ridership by route/line. Results should look something like:

```
long <- circ %>%
  pivot longer(
           cols = starts_with(c("orange", "purple", "green", "banner")),
        names_to = "bus_type",
       values_to = "number_of_individuals")
long
# A tibble: 13,752 × 5
   day date daily bus_type
                                           number of individuals
   <chr> <chr> <dbl> <chr>
                                                           <dbl>
                      952 orangeBoardings
 1 Monday 01/11/2010
                                                             877
 2 Monday 01/11/2010
                      952 orangeAlightings
                                                            1027
 3 Monday 01/11/2010
                      952 orangeAverage
                                                             952
 4 Monday 01/11/2010
                      952 purpleBoardings
                                                              NA
 5 Monday 01/11/2010
                      952 purpleAlightings
                                                              NA
 6 Monday 01/11/2010
                      952 purpleAverage
                                                              NA
 7 Monday 01/11/2010
                      952 greenBoardings
                                                              NA
                      952 greenAlightings
 8 Monday 01/11/2010
                                                              NA
 9 Monday 01/11/2010
                      952 greenAverage
                                                              NA
10 Monday 01/11/2010
                      952 bannerBoardings
                                                              NA
# ... with 13,742 more rows
```

pivot\_wider...

#### Reshaping data from long to wide

pivot\_wider() - spreads row data into columns (tidyr package)

- names\_from = the old column whose contents will be spread into multiple new column names.
- values\_from = the old column whose contents will fill in the values of those new columns.

#### Reshaping data from long to wide

0.516 0.514

1

```
long_data
# A tibble: 3 \times 2
 Month
                 Rate
       <dbl>
 <chr>
1 June_vacc_rate 0.516
2 May_vacc_rate 0.514
3 April_vacc_rate 0.511
wide_data <- long_data %>% pivot_wider(names_from = "Month",
                                    values_from = "Rate")
wide_data
# A tibble: 1 \times 3
  June_vacc_rate May_vacc_rate April_vacc_rate
          <dbl>
```

0.511

#### Reshaping Charm City Circulator

#### long

```
# A tibble: 13,752 × 5
   dav
          date daily bus_type
                                            number_of_individuals
   <chr> <chr> <dbl> <chr>
                                                            <dbl>
 1 Monday 01/11/2010
                      952 orangeBoardings
                                                              877
 2 Monday 01/11/2010
                      952 orangeAlightings
                                                             1027
 3 Monday 01/11/2010
                      952 orangeAverage
                                                              952
 4 Monday 01/11/2010
                       952 purpleBoardings
                                                               NA
 5 Monday 01/11/2010
                       952 purpleAlightings
                                                               NA
                       952 purpleAverage
 6 Monday 01/11/2010
                                                               NA
 7 Monday 01/11/2010
                       952 greenBoardings
                                                               NA
 8 Monday 01/11/2010
                       952 greenAlightings
                                                               NA
 9 Monday 01/11/2010
                       952 greenAverage
                                                               NA
10 Monday 01/11/2010
                       952 bannerBoardings
                                                               NA
# ... with 13,742 more rows
```

#### Reshaping Charm City Circulator

```
wide <- long %>% pivot_wider(names_from = "bus_type",
                            values_from = "number_of_individuals")
wide
# A tibble: 1,146 × 15
            date
                       daily orangeBoardings orangeAlightings orangeAverage
   dav
   <chr>
            <chr>
                       <dbl>
                                       <dbl>
                                                        <dbl>
                                                                      <dbl>
 1 Monday
            01/11/2010 952
                                         877
                                                         1027
                                                                       952
 2 Tuesday
            01/12/2010 796
                                         777
                                                         815
                                                                       796
 3 Wednesday 01/13/2010 1212.
                                        1203
                                                         1220
                                                                      1212.
 4 Thursday
            01/14/2010 1214.
                                        1194
                                                         1233
                                                                      1214.
 5 Friday
            01/15/2010 1644
                                        1645
                                                         1643
                                                                      1644
 6 Saturday
            01/16/2010 1490.
                                        1457
                                                         1524
                                                                      1490.
 7 Sunday
                        888.
                                       839
                                                                      888.
            01/17/2010
                                                         938
 8 Monday
            01/18/2010 1000.
                                       999
                                                         1000
                                                                      1000.
 9 Tuesday
            01/19/2010 1035
                                        1023
                                                         1047
                                                                      1035
10 Wednesday 01/20/2010 1396.
                                        1375
                                                         1416
                                                                      1396.
# ... with 1,136 more rows, and 9 more variables: purpleBoardings <dbl>,
    purpleAlightings <dbl>, purpleAverage <dbl>, greenBoardings <dbl>,
#
    greenAlightings <dbl>, greenAverage <dbl>, bannerBoardings <dbl>,
#
#
    bannerAlightings <dbl>, bannerAverage <dbl>
```

#### Summary

- pivot\_longer() goes from wide -> long
  - Specify columns you want to pivot
  - Specify names\_to = and values\_to = for custom naming
- pivot\_wider() goes from long -> wide
  - Specify names\_from = and values\_from =

Workshop Website