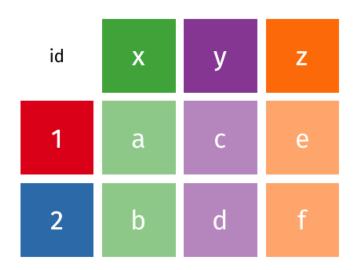
# Data Reshaping

Data Wrangling in R

## Reshaping: wide vs. long data

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

wide



#### What is wide/long data?

Data is stored *differently* in the tibble.

Wide: has many columns

#### Long: column names become data

#### What is wide/long data?

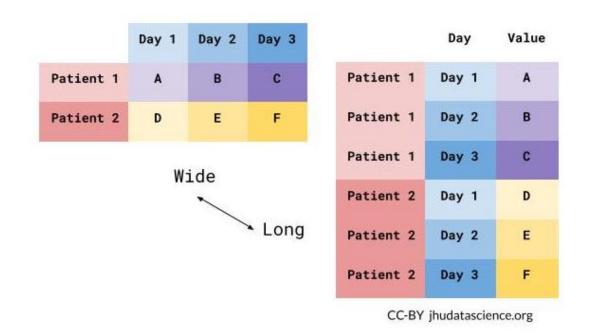
Wide: multiple columns per individual, values spread across multiple columns

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

```
# A tibble: 6 × 3
State name value
<chr> <chr> <chr> 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
```

#### What is wide/long data?

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



#### 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 × 3
State name value
<chr> <chr> 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\_wider make a variable into multiple columns, (long to wide)
- separate string into multiple columns

The reshape command exists. Its arguments are considered more confusing, so we don't recommend it.

You might see old functions gather and spread when googling. These are older iterations of pivot longer and pivot wider, respectively.

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})
```

```
wide data
\# A tibble: 1 \times 3
  June_vacc_rate May_vacc_rate April_vacc_rate

    <dbl>

    0.516

    <dbl>

    0.514

                                            \overline{<}dbl>
                                              0.511
long data <- wide data %>% pivot longer(cols = everything())
long data
\# A tibble: 3 \times 2
  name value
 <chr> <dbl>
1 June vacc rate 0.516
2 May_vacc rate 0.514
3 April_vacc_rate 0.511
```

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

```
wide data
\# A tibble: 1 \times 3
 June_vacc_rate May_vacc_rate April_vacc_rate
          \overline{<}dbl>
          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
2 May vacc rate 0.514
3 April vacc_rate 0.511
```

Newly created column names are enclosed in quotation marks.

#### Data used: Charm City Circulator

https://sisbid.github.io/Data-Wrangling/data/Charm\_City\_Circulator\_Ridership.csv

```
circ <-
  read csv("https://sisbid.github.io/Data-Wrangling/data/Charm City Circulator Ridership.csv")
head(circ, 5)
\# A tibble: 5 \times 15
           date orangeBoardings orangeAlightings orangeAverage purpleBoardings purpleAlightings purpleAverage
  day
  <chr>
                                              <dbl>
                                                             <dbl>
           <chr>
                            <dbl>
                                                                             <dbl>
                                                                                               <dbl>
                                                                                                              <dbl>
1 Monday 01/1...
                              877
                                               1027
                                                              952
                                                                                NA
                                                                                                  NA
                                                                                                                 NA
2 Tuesday 01/1...
                              777
                                                815
                                                             796
                                                                                NA
                                                                                                  NA
                                                                                                                 NA
3 Wednesd... 01/1...
                             1203
                                               1220
                                                            1212.
                                                                                NA
                                                                                                  NA
                                                                                                                 NA
4 Thursday 01/1...
                                                            1214.
                             1194
                                               1233
                                                                                NA
                                                                                                  NA
                                                                                                                 NA
5 Friday 01/1...
                                               1643
                                                            1644
                             1645
                                                                                NA
                                                                                                  NA
                                                                                                                 NA
# i 7 more variables: greenBoardings <dbl>, greenAlightings <dbl>, greenAverage <dbl>, bannerBoardings <dbl>,
    bannerAlightings <dbl>, bannerAverage <dbl>, daily <dbl>
```

```
long <- circ %>%
 pivot longer(starts with(c("orange", "purple", "green", "banner")))
long
# A tibble: 13,752 \times 5
  day date daily name
                                          value
  <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
 6 Monday 01/11/2010
                      952 purpleAverage
                                             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
# i 13,742 more rows
```

There are many ways to select the columns we want. Use <code>?tidyr\_tidy\_select</code> to look at more column selection options.

```
long <- circ %>%
 pivot longer (!c(day, date, daily))
long
# A tibble: 13,752 \times 5
  day date daily name
                                           value
  <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
 6 Monday 01/11/2010
                       952 purpleAverage
                                              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
# i 13,742 more rows
```

#### Cleaning up long data

We will use str replace from the stringr package to put in the names

```
long <- long %>% mutate(
 name = str replace(name, "Board", " Board"),
 name = str replace(name, "Alight", " Alight"),
 name = str replace(name, "Average", " Average")
long
# A tibble: 13,752 × 5
  day date daily name
                                           value
  <chr> <chr> <chr> <dbl> <chr>
                                          <dbl>
1 Monday 01/11/2010 952 orange Boardings 877
2 Monday 01/11/2010 952 orange Alightings 1027
3 Monday 01/11/2010 952 orange Average
                                            952
 4 Monday 01/11/2010
                     952 purple Boardings
                                              NA
 5 Monday 01/11/2010
                      952 purple Alightings
                                           NA
 6 Monday 01/11/2010
                      952 purple Average
                                              NA
7 Monday 01/11/2010
                      952 green Boardings
                                              NA
8 Monday 01/11/2010
                      952 green Alightings
                                              NA
                     952 green Average
 9 Monday 01/11/2010
                                              NA
10 Monday 01/11/2010
                      952 banner Boardings
                                              NA
# i 13,742 more rows
```

#### Cleaning up long data

Now each var is Boardings, Averages, or Alightings. We use "into =" to name the new columns and "sep =" to show where the separation should happen.

```
long <- long %>%
  separate(name, into = c("line", "type"), sep = " ")
long
# A tibble: 13,752 × 6
  day date daily line type value
  <chr> <chr> <chr> <dbl> <chr> <dbl>
 1 Monday 01/11/2010 952 orange Boardings
                                            877
 2 Monday 01/11/2010 952 orange Alightings 1027
 3 Monday 01/11/2010 952 orange Average
                                          952
 4 Monday 01/11/2010
                     952 purple Boardings
                                             NA
 5 Monday 01/11/2010
                      952 purple Alightings
                                             NA
 6 Monday 01/11/2010
                      952 purple Average
                                             NA
                      952 green Boardings
 7 Monday 01/11/2010
                                             NA
 8 Monday 01/11/2010
                      952 green Alightings
                                             NA
 9 Monday 01/11/2010
                    952 green Average
                                             NA
10 Monday 01/11/2010
                     952 banner Boardings
                                             NA
# i 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

```
long data
\# A tibble: 3 \times 2
 Month Rate
 <chr> <dbl>
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
         - \overline{\langle}dbl\rangle
          0.516 0.514
                                   0.511
```

#### Reshaping Charm City Circulator

#### long

```
# A tibble: 13,752 \times 6
  day date daily line type
                                            value
  <chr> <chr> <chr> <dbl> <chr>
                                           <dbl>
1 Monday 01/11/2010
                      952 orange Boardings
                                              877
2 Monday 01/11/2010
                      952 orange Alightings
                                             1027
 3 Monday 01/11/2010
                      952 orange Average
                                              952
                      952 purple Boardings
 4 Monday 01/11/2010
                                               NA
 5 Monday 01/11/2010
                      952 purple Alightings
                                               NA
 6 Monday 01/11/2010
                      952 purple Average
                                               NA
7 Monday 01/11/2010
                      952 green Boardings
                                               NA
 8 Monday 01/11/2010
                      952 green Alightings
                                               NA
 9 Monday 01/11/2010
                      952 green Average
                                               NA
10 Monday 01/11/2010
                      952 banner Boardings
                                               NA
# i 13,742 more rows
```

#### Reshaping Charm City Circulator

```
wide <- long %>% pivot wider (names from = "type",
                           values from = "value")
wide
\# A tibble: 4,584 \times 7
            date
                      daily line Boardings Alightings Average
  day
                      <dbl> <chr>
                                       <dbl>
  <chr>
          <chr>
                                                 <dbl>
                                                         <dbl>
 1 Monday 01/11/2010 952 orange
                                         877
                                                  1027
                                                          952
 2 Monday 01/11/2010 952 purple
                                                           NA
                                         NA
                                                    NA
            01/11/2010 952 green
 3 Monday
                                         NA
                                                    NA
                                                          NA
 4 Monday
            01/11/2010 952 banner
                                         NA
                                                    NA
                                                          NA
 5 Tuesday
            01/12/2010 796 orange
                                        777
                                                   815
                                                          796
 6 Tuesday
            01/12/2010 796 purple
                                         NA
                                                    NA
                                                          NA
 7 Tuesday
            01/12/2010 796 green
                                         NA
                                                    NA
                                                          NA
            01/12/2010 796 banner
 8 Tuesday
                                         NA
                                                    NA
                                                          NA
 9 Wednesday 01/13/2010 1212. orange
                                       1203
                                                  1220
                                                         1212.
10 Wednesday 01/13/2010 1212. purple
                                         NA
                                                    NA
                                                           NA
# i 4,574 more rows
```

# Adding prefixes

## Prefixes when pivoting

the datasets::airquality data shows various air quality metrics measured in New York in 1973.

```
air <- datasets::airquality %>% select(Temp, Month, Day)
air
```

	Temp	Month	Day
1	67	5	$\bar{1}$
2	72	5	2
1 2 3 4 5 6 7 8 9 10	74	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	1 2 3 4 5 6 7
4	74 62	5	4
5	56 66 65	5	5
6	66	5	6
7	65	5	7
8	59	5	8
9	61	5	
10	69	5	10
11	74	5	11
12	69	5	12
13	66	5	13
14	68	5	14
15	58	5	15
16 17	64	5	16
17	66	5	17
18	57	5	18
19	68	5	19
20 21	62 59	5	20 21
21	59	.5	21

#### Prefixes when pivoting

Let's pivot Month wider: but it might be helpful to add "Month" to the new column name so it isn't just numbers.

```
air %>% pivot wider(names from = "Month",
                     values from = "Temp")
\# A tibble: 31 \times 6
     Day `5`
   <int> <int> <int> <int> <int><</pre>
                   78
                          84
                                81
                                       91
             67
                         85
                                81
                                      92
                74
                67
                                82
            74
                         81
                                      93
                84
                         84
                                86
                                      93
            56
                85
                         83
                                85
                                      87
            66
                   79
                         83
                                87
                                      84
                   82
            65
                         88
                                89
                                      80
       8
                   87
            59
                          92
                                90
                                      78
 9
       9
            61
                   90
                          92
                                90
                                      75
10
      10
                   87
                          89
                                      73
            69
                                92
    21 more rows
```

#### Prefixes when pivoting

#### Much better!

```
air %>% pivot wider(names from = "Month",
                values from = "Temp",
                names prefix = "Month ")
\# A tibble: 31 \times 6
    Day Month 5 Month 6 Month 7 Month 8 Month 9
  \langle int \rangle \langle int \rangle \langle int \rangle \langle int \rangle
           67
                  78
                         84
                                      91
                               81
                  74 85 81
                                      92
       74 67 81 82
                                      93
        62
              84 84 86
                                      93
        56
              85 83 85
                                  87
                79
                       83 87
         66
                                      84
                       88 89
               82
         65
                                      80
8
                        92 90
     8
         59
              87
                                      78
9
         61
                  90
                       92 90
                                      75
10
        69
                  87
                         89
                            92
                                      73
# i 21 more rows
```

#### Summary

- tidyr package helps us convert between wide and long data
- 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 =