Merging Data Sets

Data Wrangling in R

Joining in dplyr

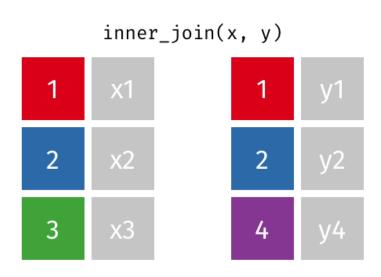
- Merging/joining data sets together usually on key variables, usually "id"
- ?join see different types of joining for dplyr
- inner_join(x, y) only rows that match for x and y are kept
- full_join(x, y) all rows of x and y are kept
- left_join(x, y) all rows of x are kept even if not merged with y
- right join(x, y) all rows of y are kept even if not merged with x
- anti join(x, y) all rows from x not in y keeping just columns from x.

Merging: Simple Data

1 Maine 32.4%
2 Alaska 41.7%

Inner Join

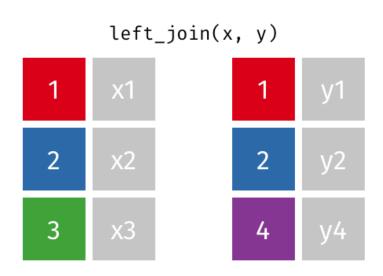
https://github.com/gadenbuie/tidyexplain/blob/main/images/inner-join.gif



Inner Join

Left Join

https://raw.githubusercontent.com/gadenbuie/tidyexplain/main/images/left-join.gif



Left Join

Install tidylog package to log outputs

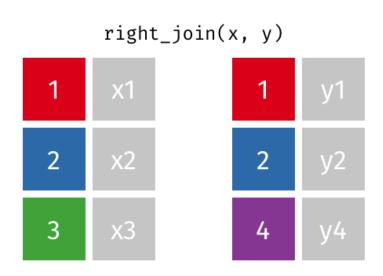
```
# install.packages("tidylog")
library(tidylog)
left join (data As, data cold)
Joining, by = "State"
left join: added one column (April vacc rate)
> rows only in x 1
> rows only in y (1)
> matched rows 1
> ===
> rows total 2
\# A tibble: 2 × 4
  State June_vacc_rate May_vacc_rate April_vacc_rate

      <chr>
      <chr>
      <chr>
      36.0%
      <NA>

2 Alaska 47.5% 46.2% 41.7%
```

Right Join

https://raw.githubusercontent.com/gadenbuie/tidyexplain/main/images/right-join.gif



Right Join

```
rj <- right join(data As, data_cold)</pre>
Joining, by = "State"
right join: added one column (April vacc rate)
> rows only in x (1)
> rows only in y 1
> matched rows 1
> ===
> rows total 2
rj
# A tibble: 2 \times 4
  State June_vacc_rate May_vacc_rate April_vacc_rate

      <chr>
      <chr>
      <chr>
      <chr>
      46.2%
      41.7%

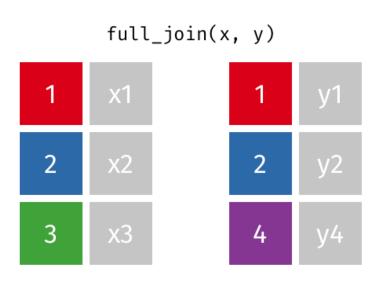
                                        32.4%
2 Maine <NA> <NA>
```

Left Join: Switching arguments

```
lj2 <- left join(data cold, data As)
Joining, by = "State"
left join: added 2 columns (June vacc rate, May vacc rate)
> rows only in x 1
> rows only in y (1)
> matched rows 1
> ===
> rows total 2
1j2
# A tibble: 2 \times 4
 State April_vacc_rate June_vacc_rate May_vacc_rate
<NA>
                               <NA>
1 Maine 32.4%
2 Alaska 41.7% 47.5% 46.2%
```

Full Join

https://raw.githubusercontent.com/gadenbuie/tidyexplain/main/images/full-join.gif



Full Join

```
fj <- full join(data As, data cold)</pre>
Joining, by = "State"
full join: added one column (April vacc rate)
> rows only in x 1
> rows only in y 1
> matched rows 1
> ===
> rows total 3
fj
# A tibble: 3 \times 4
 State June_vacc_rate May_vacc_rate April_vacc_rate
                  \langle chr \rangle \langle chr \overline{\rangle}
1 Alabama 37.2%
                     46.2% 41.7%
2 Alaska 47.5%
3 Maine <NA> <NA> 32.4%
```


data As

```
lj <- left_join(data_As, data_cold)

Joining, by = "State"
left_join: added 2 columns (vacc_rate, month)
> rows only in x 1
> rows only in y (1)
> matched rows 2 (includes duplicates)
> ===
> rows total 3
```

Data including the joining column ("State") has been duplicated.

1 j

Note that "Alaska willow ptarmigan" appears twice.

https://github.com/gadenbuie/tidyexplain/blob/main/images/left-join-extra.gif

<pre>left_join(x, y)</pre>			
1	x1	1	y1
2	x2	2	у2
3	х3	4	у4
		2	у5

Stop tidylog

unloadNamespace("tidylog")

Using the by argument

By default joins use the intersection of column names. If by is specified, it uses that.

Using the by argument

You can join based on multiple columns by using something like by = c(col1, col2).

If the datasets have two different names for the same data, use:

```
full_join(x, y, by = c("a" = "b"))
```

Using "setdiff"

We might want to determine what indexes ARE in the first dataset that AREN'T in the second:

```
data_As

# A tibble: 2 × 2
   State   state_bird
   <chr>   <chr>
1 Alabama wild turkey
2 Alaska willow ptarmigan

data_cold

# A tibble: 3 × 3
   State vacc_rate month
   <chr>   <chr>   <chr>
1 Maine 32.4%   April
2 Alaska 41.7%   April
3 Alaska 46.2%   May
```

Using "setdiff"

Use setdiff to determine what indexes ARE in the first dataset that AREN'T in the second:

```
A_states <- data_As %>% pull(State)
cold_states <- data_cold %>% pull(State)

setdiff(A_states, cold_states)

[1] "Alabama"

setdiff(cold_states, A_states)

[1] "Maine"
```

Using bind_rows()(dplyr)

```
one <- starwars[1:4, ]
two <- starwars[9:12, ]
bind rows (one, two)
# A tibble: 8 × 14
  name height mass hair color skin color eye color birth year sex
                                                                          geno
  \langle chr \rangle \langle int \rangle \langle chr \rangle \langle chr \rangle \langle chr \rangle
                                                     <dbl> <chr> <chr
1 Luke Sky... 172 77 blond fair blue
2 C-3PO 167 75 <NA> gold yellow
3 R2-D2 96 32 <NA> white, bl... red
                                                             19
                                                                    male
                                                                          maso
                                                          112
                                                                   none
                                                                          maso
                                                               33
                                                                    none
                                                                          maso
4 Darth Va... 202 136 none white yellow 41.9 male
                                                                          maso
5 Biggs Da... 183 84 black light brown
                                                               24
                                                                    male
                                                                          maso
6 Obi-Wan ... 182 77 auburn, w... fair blue-gray
                                                               57
                                                                    male
                                                                          maso
7 Anakin S... 188 84 blond
                                    fair blue
                                                           41.9 male
                                                                          maso
8 Wilhuff ... 180 NA auburn, g... fair blue
                                                               64
                                                                    male
                                                                          maso
# ... with 5 more variables: homeworld <chr>, species <chr>, films <list>,
    vehicles <list>, starships <list>
```

Summary

- Merging/joining data sets together assumes all column names that overlap
 - use the by = c("a" = "b") if they differ
- inner_join(x, y) only rows that match for x and y are kept
- full_join(x, y) all rows of x and y are kept
- · left_join(x, y) all rows of x are kept even if not merged with y
- right_join(x, y) all rows of y are kept even if not merged with x
- Use the tidylog package for a detailed summary
- setdiff(x, y) shows what in x is missing from y
- bind rows(x, y) appends datasets