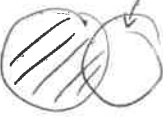# 7: Combining Data Sets with `dplyr`

Suppose you have the following two data sets. The first, `df1` has the variables `id_numb` and `xvar`. The second, `df2` has the variables `id` and `yvar`. `id_numb` and `id` serve as identification variables, possibly with duplicates, where observations from the first data set with `id_numb = 1` correspond to observations in the second data set with `id = 1`.

| id_numb | xvar |
|---------|------|
| 1 | 16 |
| 1 | -1 |
| 2 | 11 |
| 4 | 13 |

| id | yvar |
|----|------|
| 1 | -1 |
| 2 | -4 |
| 2 | 0 |
| 3 | -9 |

**Mutating Joins**

- `left_join()`

dropped

$$\text{left\_join}(df1, df2), by = c(\underbrace{"id\_numb"}_{key\ in\ df1} = \underbrace{"id"}_{key\ in\ df2}))$$

$\Rightarrow$

| id_numb | xvar | yvar |
|---------|------|------|
| 1 | 16 | -1 |
| 1 | -1 | -1 |
| 2 | 11 | -4 |
| 2 | 11 | 0 |
| 4 | 13 | ~~0~~ NA |

- `right_join()` ( id=3 dropped)

dropped right_join(df1, df2)

$\Longleftrightarrow$ left_join (df2, df1)

- `inner_join()`

inner_join(df1, df2, by = c("id_numb" = "id"))

$\Rightarrow$

1

| id_numb | xvar |
|---------|------|
| 1 | 16 |
| 1 | -1 |
| 2 | 11 |
| 4 | 13 |

| id | yvar |
|----|------|
| 1 | -1 |
| 2 | -4 |
| 2 | 0 |
| 3 | -9 |

- `full_join()`



Same as left_join() but adds

| id_numb | xvar | yvar |
|---------|------|------|
| ~~3~~ 3 | NA | -9 |

**Filtering Joins**

- `semi_join()` ⟋ same syntax     keep anything that has a match in df2

· df1, df2 =>

| id_numb | xvar |
|---------|------|
| 1 | 16 |
| 1 | -1 |
| 2 | 11 |

- `anti_join()` ⟋ same syntax     keep anything that doesn't have a match in df2

df1, df2 =>

| id_numb | xvar |
|---------|------|
| 4 | 13 |