

Merging datasets

Many datasets are split across different tables/files and need to be combined during analysis.

The `join` family of functions provides all the functionality you need.

When merging two tables, first decide what columns define a **match** between rows.

```
blah_join(x, y, by = <columns>)
```

Pick from six join types, which differ in what they do with **non-matching** rows.

Join type	Rows of x	Columns of x	Columns of y
inner	> 1 match in y	all	all
left	all	all	all
right	> 1 match in y	all	all
full	all	all	all
semi	> 1 match in y	all	none
anti	0 match in y	all	none

sightings		pokedex	
location	poke_id	ID	name
ORC	2	1	Bulbasaur
ORC	1	2	Squirtle
Sloan	1	3	Charmander
Sloan	4	?	?

We would like to match `poke_id` on the LHS to `ID` on the RHS, so we specify:

```
xxxx_join(sightings, pokedex, by = c("poke_id" = "ID"))
```



sightings		pokedex	
location	poke_id	ID	name
ORC	2	1	Bulbasaur
ORC	1	2	Squirtle
Sloan	1	3	Charmander
Sloan	4	?	?

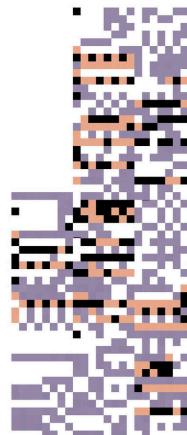
```
sightings %>%
  inner_join(pokedex, by = c("poke_id" = "ID"))
```

location	poke_id	name
ORC	2	Squirtle
ORC	1	Bulbasaur
Sloan	1	Bulbasaur

sightings		pokedex	
location	poke_id	ID	name
ORC	2	1	Bulbasaur
ORC	1	2	Squirtle
Sloan	1	3	Charmander
Sloan	4	?	?

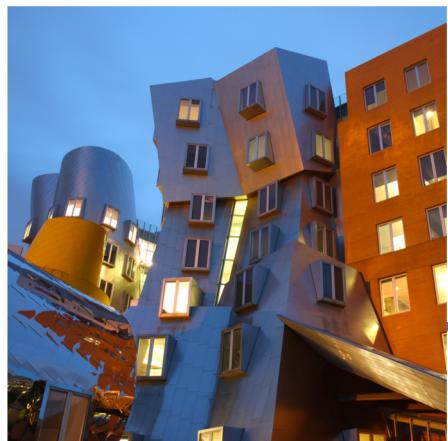
```
sightings %>%
  left_join(pokedex, by = c("poke_id" = "ID"))
```

location	poke_id	name
ORC	2	Squirtle
ORC	1	Bulbasaur
Sloan	1	Bulbasaur
Sloan	4	NA



sightings		pokedex	
location	poke_id	ID	name
ORC	2	1	Bulbasaur
ORC	1	2	Squirtle
Sloan	1	3	Charmander
Sloan	4	?	?

```
sightings %>%
  right_join(pokedex, by = c("poke_id" = "ID"))
```



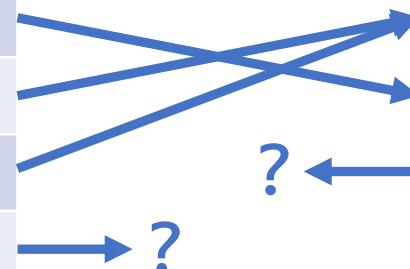
location	poke_id	name
ORC	2	Squirtle
ORC	1	Bulbasaur
Sloan	1	Bulbasaur
NA	3	Charmander

sightings

location	poke_id
ORC	2
ORC	1
Sloan	1
Sloan	4

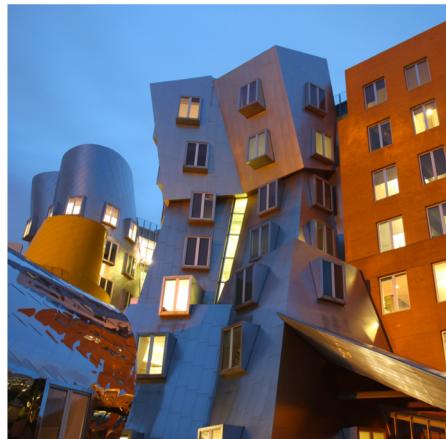
pokedex

ID	name
1	Bulbasaur
2	Squirtle
3	Charmander

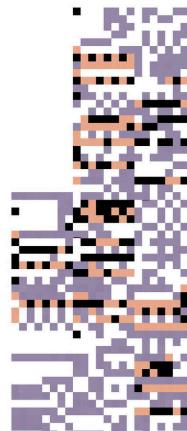


sightings %>%

full_join(pokedex, by = c("poke_id" = "ID"))



location	poke_id	name
ORC	2	Squirtle
ORC	1	Bulbasaur
Sloan	1	Bulbasaur
Sloan	4	NA
NA	3	Charmander



sightings		pokedex	
location	poke_id	ID	name
ORC	2	1	Bulbasaur
ORC	1	2	Squirtle
Sloan	1	3	Charmander
Sloan	4	? → ?	

```
sightings %>%
  semi_join(pokedex, by = c("poke_id" = "ID"))
```

location	poke_id
ORC	2
ORC	1
Sloan	1

sightings		pokedex	
location	poke_id	ID	name
ORC	2	1	Bulbasaur
ORC	1	2	Squirtle
Sloan	1	3	Charmander
Sloan	4	? → ?	

```
sightings %>%
  anti_join(pokedex, by = c("poke_id" = "ID"))
```

location	poke_id
Sloan	4

sightings		types	
location	name	name	type
ORC	Squirtle	Bulbasaur	Grass
ORC	Bulbasaur	Bulbasaur	Poison
Sloan	Bulbasaur	Squirtle	Water
Sloan	Pikachu	Charmander	Fire

```

graph LR
    S1[location: ORC, name: Squirtle] --> T1["name: Bulbasaur, type: Grass"]
    S1 --> T2["name: Bulbasaur, type: Poison"]
    S2[location: ORC, name: Bulbasaur] --> T1
    S2 --> T3["name: Squirtle, type: Water"]
    S3[Sloan, name: Bulbasaur] --> T1
    S3 --> T3
    S4[Sloan, name: Pikachu] --> T4["name: Charmander, type: Fire"]
  
```

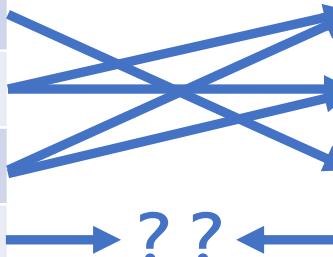
If there are **multiple** matches to a row, the result will have a row for **every combination**.

sightings

location	name
ORC	Squirtle
ORC	Bulbasaur
Sloan	Bulbasaur
Sloan	Pikachu

types

name	type
Bulbasaur	Grass
Bulbasaur	Poison
Squirtle	Water
Charmander	Fire



```
sightings %>%  
  inner_join(types, by = "name")
```

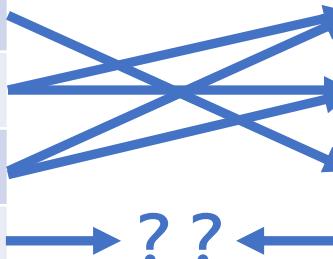
location	name	type
ORC	Squirtle	Water
ORC	Bulbasaur	Grass
ORC	Bulbasaur	Poison
Sloan	Bulbasaur	Grass
Sloan	Bulbasaur	Poison

sightings

location	name
ORC	Squirtle
ORC	Bulbasaur
Sloan	Bulbasaur
Sloan	Pikachu

types

name	type
Bulbasaur	Grass
Bulbasaur	Poison
Squirtle	Water
Charmander	Fire



sightings %>%

```
left_join(types, by = "name")
```

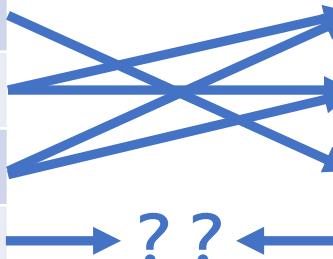
location	name	type
ORC	Squirtle	Water
ORC	Bulbasaur	Grass
ORC	Bulbasaur	Poison
Sloan	Bulbasaur	Grass
Sloan	Bulbasaur	Poison
Sloan	Pikachu	NA

sightings

location	name
ORC	Squirtle
ORC	Bulbasaur
Sloan	Bulbasaur
Sloan	Pikachu

types

name	type
Bulbasaur	Grass
Bulbasaur	Poison
Squirtle	Water
Charmander	Fire



sightings %>%

inner_join(types, by = "name")

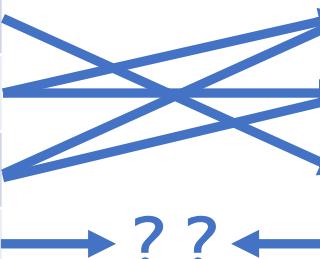
location	name	type
ORC	Squirtle	Water
ORC	Bulbasaur	Grass
ORC	Bulbasaur	Poison
Sloan	Bulbasaur	Grass
Sloan	Bulbasaur	Poison

%>% count(type)

type	n
Water	1
Grass	2
Poison	2

sightings

location	name
ORC	Squirtle
ORC	Bulbasaur
Sloan	Bulbasaur
Sloan	Pikachu



types

name	type
Bulbasaur	Grass
Bulbasaur	Poison
Squirtle	Water
Charmander	Fire

sightings %>%

```
inner_join(types, by = "name")
```

```
%>% count(name)
```

location	name	type
ORC	Squirtle	Water
ORC	Bulbasaur	Grass
ORC	Bulbasaur	Poison
Sloan	Bulbasaur	Grass
Sloan	Bulbasaur	Poison

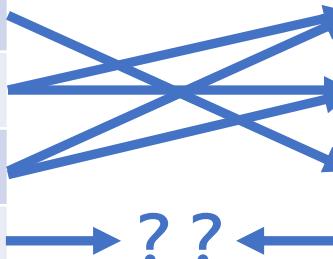
type	n
Squirtle	1
Bulbasaur	4

sightings

location	name
ORC	Squirtle
ORC	Bulbasaur
Sloan	Bulbasaur
Sloan	Pikachu

types

name	type
Bulbasaur	Grass
Bulbasaur	Poison
Squirtle	Water
Charmander	Fire



```
sightings %>%
  inner_join(types, by = "name")
```

location	name	type
ORC	Squirtle	Water
ORC	Bulbasaur	Grass
ORC	Bulbasaur	Poison
Sloan	Bulbasaur	Grass
Sloan	Bulbasaur	Poison

Moral

Always be aware of row duplication!!!

sightings		catch			
location	name		name	location	catch_pct
ORC	Squirtle		Bulbasaur	ORC	50
ORC	Bulbasaur		Bulbasaur	Sloan	75
Sloan	Bulbasaur		Squirtle	ORC	90
Sloan	Pikachu	→ ?	Squirtle	Sloan	100
			Pikachu	ORC	42

A match can be defined by more than one column.

```
sightings %>%
  left_join(catch, by = c("name" = "name", "location" = "location"))
```

location	name	catch_pct
ORC	Squirtle	90
ORC	Bulbasaur	50
Sloan	Bulbasaur	75
Sloan	Pikachu	NA