## **Practicing and Extending Joins**DATA 202 21FA

#### Reminders

Data Science Firepit Tonight!

#### Also:

- Quiz 5 closes tonight
- Homework 3 due Friday
- Discussion Forum replies due tomorrow

## Preview: replication project

## Gradescope trial

See Moodle.



Can't you rewrite any right\_join as a flipped left\_join?

Yes! It's mostly convenience (piped data comes in on *left*).

Can we join more than 2 tables at once?

SQL? yes. dplyr makes you split it up though.

# Q: Why did we need summarize()?

[1] 5.75

```
purchases_and_prices %>%
  summarize(total_revenue = sum(price))
# A tibble: 1 \times 1
  total revenue
          <dbl>
         5.75
sum(price)
Error in eval(expr, envir, enclos): object 'price' not found
sum(purchases_and_prices$price)
```

#### **Aside: vectors**

So far we've always packaged vectors in data frames. But they can live on their own.

```
price_vec<- purchases_and_pric
price_vec</pre>
```

```
[1] 1.00 0.80 0.15 0.80 3.00
```

```
typeof(price_vec)
```

```
[1] "double"
```

#### **Relational Data**

- Bring together data that may never have been brought together before
- Values:
  - No one source has a complete view of the world
  - We need each other ("one body, many parts")
- Helps us get a more holistic picture
  - people aren't just their transactions; they also live somewhere

#### **Example Applications**

- Looking up abbreviations
- Connecting sales data to each customer's demographics
- Combining student data from Moodle, an online textbook,
   GitHub activity reports, ...

What others can you think of?

#### **Code Together: Flight Delays**

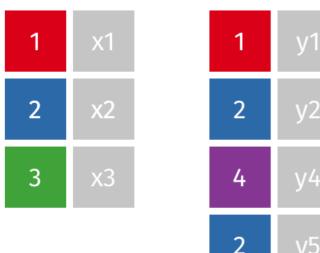
#### library(nycflights13)

```
flights %>%
  drop_na(arr_delay) %>%
  group_by(carrier) %>%
  summarize(avg_delay = mean(arr_delay)) %>%
  arrange(desc(avg_delay)) %>%
  left_join(airlines, by = "carrier") %>%
  select(name, avg_delay)
```

```
# A tibble: 16 \times 2
                               avg_delay
  name
  <chr>
                                    <dbl>
1 Frontier Airlines Inc.
                                     21.9
2 AirTran Airways Corporation
                                    20.1
3 ExpressJet Airlines Inc.
                                15.8
4 Mesa Airlines Inc.
                                    15.6
5 SkyWest Airlines Inc.
                                    11.9
6 Envoy Air
                                     10.8
# ... with 10 more rows
```

## Multiple matches

left\_join(x, y)



## Revenue by item?

purchases

customer_id	item
c1	bread
c1	milk
c1	banana
c2	milk
c2	toilet paper

prices

item	price
avocado	0.50
banana	0.15
bread	1.00
milk	0.80
toilet paper	3.00

#### Revenue by item?

prices %>%

For each item, look up all sales data.

Notice: multiple rows for each item. Where did each one come from?

#### Revenue by item?

```
prices %>%
  left_join(purchases) %>%
  group_by(item) %>%
  summarize(revenue = sum(price))
```

...but nobody bought any avocados!

#### inner vs left

```
prices %>%
  left_join(purchases)
```

```
prices %>%
  inner_join(purchases)
```

#### Put it together

```
prices %>%
  inner_join(purchases) %>%
  group_by(item) %>%
  summarize(revenue = sum(price))
```

## uh oh...

customer_id	item
c1	bread
c1	milk
c1	bananas
c2	milk
c2	toilet paper

item	price
AVOCADO	0.50
BANANA	0.15
BREAD	1.00
MILK	0.80
TOILET_PAPER	3.00

#### uh oh...

customer_id	item
c1	bread
c1	milk
c1	bananas
c2	milk
c2	toilet paper

item	price
AVOCADO	0.50
BANANA	0.15
BREAD	1.00
MILK	0.80
TOILET_PAPER	3.00

```
purchases %>%
  left_join(prices)
```

```
# A tibble: 5 \times 3
  customer_id item
                           price
         <chr>
  <chr>
                           <dbl>
            bread
1 c1
                              NA
              milk
2 c1
                              NA
3 c1
              bananas
                              NA
              milk
4 c2
                              NA
              toilet paper
5 c2
                              NA
```

#### **Specifying keys**

- Keys must match exactly
- Can join on multiple columns (first name and last name)
- Default join: columns with same names
- Specify what columns to use: left\_join(x, y, by =
  c("first\_name", "last\_name"))

#### Make a new key that does match

```
# A tibble: 5 \times 4
                      item.y price
 customer id item
 <chr> <chr>
                      <chr> <dbl>
        bread
1 c1
                      BREAD
                             1
2 c1 milk
                     MILK 0.8
3 c1
           bananas <NA>
                            NA
           milk
                     MILK 0.8
4 c2
5 c2
           toilet paper <NA>
                            NA
```

#### But still some mismatches

```
purchases %>%
   anti_join(
   prices_fixed,
   by = c("item" = "item_norm")
```

2 BANANA 0.15 banana

3 TOILET\_PAPER 3 toilet\_par

## Manual massaging: if\_else

```
prices_fixed <- prices %>%
  mutate(
    item_norm = str_to_lower(item),
    item_norm = if_else(
      item_norm == "toilet_paper", "toilet paper",
                                 item norm)
purchases %>%
  left_join(prices_fixed, by = c("item" = "item_norm"))
# A tibble: 5 \times 4
 customer id item
                         item.y price
 <chr>
        <chr>
                         <chr>
                                   <dbl>
1 c1
          bread
                         BREAD
2 c1
            milk
                        MILK
                                      0.8
3 c1
            bananas <NA>
                                      NA
4 c2
            milk
                                   0.8
                         MILK
5 c2
            toilet paper TOILET PAPER
                                       3
```

## Manual massaging: case\_when

```
prices_fixed <- prices %>%
  mutate(
    item_norm = str_to_lower(item),
    item_norm = case_when(
      item_norm == "toilet_paper" ~ "toilet paper",
     TRUE
                           ~ item norm)
purchases %>%
  left_join(prices_fixed, by = c("item" = "item_norm"))
# A tibble: 5 \times 4
 customer id item
                       item.y price
 <chr>
       <chr>
                       <chr>
                                 <dbl>
1 c1
        bread
                       BREAD
2 c1
      milk
                       MILK
                                   0.8
3 c1
           bananas <NA>
                                    NA
4 c2
            milk
                                 0.8
                       MILK
5 c2
        toilet paper TOILET_PAPER
                                     3
```

#### Easy to add cases

```
prices_fixed <- prices %>%
  mutate(
    item_norm = str_to_lower(item),
    item norm = case when(
     item_norm == "toilet_paper" ~ "toilet paper",
     item norm == "banana" ~ "bananas",
                             ~ item norm)
     TRUE
purchases %>%
  left_join(prices_fixed, by = c("item" = "item_norm"))
# A tibble: 5 \times 4
 customer id item
                       item.y price
 <chr> <chr>
                       <chr> <dbl>
1 c1 bread
                       BREAD
2 c1 milk
                      MILK
                                  0.8
3 c1
           bananas
                                 0.15
                      BANANA
     milk
4 c2
                                  0.8
                      MILK
5 c2
        toilet paper TOILET PAPER 3
```

# Aside: Could you use a join for this instead?

Think about this on your own.

#### case\_when

if-elif version (Python):

```
if age < 0:
    return "invalid"
elif age < 18:
    return "child"
else:
    return "adult"</pre>
```

case\_when version:

[1] "adult"

- first to True wins in both versions
- TRUE corresponds to else (the default)

#### case\_when vectorizes

Like many R functions, it actually applies to all elements of a vector.

```
age <- c(-1, 0, 17, 18) # a vector
case_when(
  age < 0 ~ "invalid",
  age < 18 ~ "child",
  TRUE ~ "adult"
)</pre>
```

```
[1] "invalid" "child" "child" "adult"
```

#### case\_when vs if\_else

You can write the same thing either way. Which do you prefer?

```
if_else:
```

```
if_else(
  age < 0, "invald",
  if_else(
    age < 18, "child",
    "other"))</pre>
```

[1] "invald" "child" "child"

case\_when:

[1] "invalid" "child" "child'

#### case\_when in a data frame

```
people <- tribble(</pre>
 ~name, ~age,
 "Allen Linford", -1,
 "Seb Dodds", 0,
 "Charleen Lockwood", 17,
 "Ridley Burgin", 18,
people %>% mutate(
 adult = case_when(
   age < 0 ~ "invalid",
    age < 18 ~ "child",
   TRUE ~ "adult"
```

#### The recoding pattern

#### More case\_when tricks

See ?case\_when for how to:

- Deal with inconsistent data types
- Efficiently encode complicated conditionals
- Reuse case\_when expressions by making a function

and more!