Working with a single data frame

Data Science in a Box datasciencebox.org



We...

have a single data frame

Want to slice it, and dice it, and juice it, and process it

Data: Hotel bookings

- Data from two hotels: one resort and one city hotel
- Observations: Each row represents a hotel booking

hotels <- read_csv("data/hotels.csv")</pre>

select, arrange, and slice



select to keep variables

```
hotels %>%
select(hotel, lead_time)
```

```
## # A tibble: 119,390 × 2
                  lead time
##
    hotel
##
     <chr>
                      <dbl>
                        342
## 1 Resort Hotel
## 2 Resort Hotel
                        737
## 3 Resort Hotel
## 4 Resort Hotel
                         13
                        14
## 5 Resort Hotel
## 6 Resort Hotel
                         14
## # ... with 119,384 more rows
```

select to exclude variables

```
hotels %>%
  select(-agent)
```

```
## # A tibble: 119,390 × 31
             is canceled lead time arrival date ye… arrival date mo…
     hotel
     <chr>
                   <fdb>>
                              <dbl>
                                                <dbl> <chr>
## 1 Resort...
                                342
                                                 2015 July
## 2 Resort...
                                737
                                                 2015 July
## 3 Resort...
                                                 2015 July
                                                 2015 July
## 4 Resort...
                                                 2015 July
## 5 Resort...
                                 14
                                 14
## 6 Resort...
                                                 2015 July
## # ... with 119,384 more rows, and 26 more variables:
       arrival date week number <dbl>,
       arrival date day of month <dbl>,
## #
       stays in weekend nights <dbl>, stays in week nights <dbl>,
## #
       adults <dbl>, children <dbl>, babies <dbl>, meal <chr>,
## #
## #
       country <chr>, market segment <chr>,
## #
       distribution channel <chr>, is repeated guest <dbl>, ...
```

select a range of variables

select(hotel:arrival date month)

```
## # A tibble: 119,390 × 5
##
     hotel is canceled lead_time arrival_date_ye... arrival_date_mo...
##
     <chr>
                    <dbl>
                               <dbl>
                                                 <dbl> <chr>
## 1 Resort...
                                 342
                                                   2015 July
## 2 Resort...
                                 737
                                                   2015 July
                                                   2015 July
## 3 Resort...
## 4 Resort...
                                                  2015 July
                                  13
                                                  2015 July
## 5 Resort...
                                  14
## 6 Resort...
                                  14
                                                  2015 July
## # ... with 119,384 more rows
```

hotels %>%

select variables with certain characteristics

```
hotels %>%
  select(starts_with("arrival"))
```

```
## # A tibble: 119,390 × 4
##
     arrival date year arrival date month arrival date week number
                 <dbl> <chr>
##
                                                                <dbl>
                  2015 July
## 1
                                                                   27
                  2015 July
## 2
                                                                   27
## 3
                  2015 July
                                                                   27
## 4
                  2015 July
                                                                   27
## 5
                  2015 July
                                                                   27
                   2015 July
## 6
                                                                   27
     ... with 119,384 more rows, and 1 more variable:
       arrival_date_day_of_month <dbl>
```



select variables with certain characteristics

```
hotels %>%
  select(ends_with("type"))
```

```
## # A tibble: 119,390 × 4
##
     reserved_room_type assigned_room_ty... deposit_type customer_type
##
    <chr>
                        <chr>
                                                      <chr>
                                         <chr>
                                         No Deposit Transient
                                         No Deposit Transient
## 2 C
## 3 A
                                         No Deposit
                                                     Transient
                                         No Deposit Transient
                                         No Deposit Transient
## 5 A
                                         No Deposit
## 6 A
                                                     Transient
## # ... with 119,384 more rows
```

Select helpers

- starts_with(): Starts with a prefix
- ends_with(): Ends with a suffix
- contains(): Contains a literal string
- num_range(): Matches a numerical range like x01, x02, x03
- one_of(): Matches variable names in a character vector
- everything(): Matches all variables
- last_col(): Select last variable, possibly with an offset
- matches (): Matches a regular expression (a sequence of symbols/characters expressing a string/pattern to be searched for within text)

See help for any of these functions for more info, e.g. ?everything.



arrange in ascending / descending order

```
hotels %>%
  select(adults, children, babies) %>%
  arrange(babies)
```

```
hotels %>%
  select(adults, children, babies) %>%
  arrange(desc(babies))
```

slice for certain row numbers

```
# first five
hotels %>%
    slice(1:5)
```

```
## # A tibble: 5 × 32
    hotel is_canceled lead_time arrival_date_ye... arrival_date_mo...
    <chr>
                  <dbl>
                            <dbl>
                                             <dbl> <chr>
## 1 Resort...
                      0 342
                                              2015 July
                         737
## 2 Resort...
                                              2015 July
## 3 Resort...
                                              2015 July
                               13
## 4 Resort...
                                              2015 July
## 5 Resort...
                               14
                                              2015 July
## # ... with 27 more variables: arrival_date_week_number <dbl>,
    arrival_date_day_of_month <dbl>,
## #
## # stays in weekend nights <dbl>, stays in week nights <dbl>,
      adults <dbl>, children <dbl>, babies <dbl>, meal <chr>,
## #
## #
      country <chr>, market_segment <chr>,
      distribution_channel <chr>, is_repeated_guest <dbl>,
## #
      previous cancellations <dbl>, ...
## #
```

In R, you can use the # for adding comments to your code. Any text following # will be printed as is, and won't be run as R code. This is useful for leaving comments in your code and for temporarily disabling certain lines of code while debugging.

```
hotels %>%
  # slice the first five rows # this line is a comment
                                 # this one doesn't run
  #select(hotel) %>%
   slice(1:5)
                                 # this line runs
## # A tibble: 5 × 32
             is canceled lead time arrival date ye... arrival date mo...
     <chr>
                   <dbl>
                              <dbl>
                                                <dbl> <chr>
## 1 Resort...
                                342
                                                 2015 July
                                737
                                                 2015 July
## 2 Resort...
                                                 2015 July
## 3 Resort...
                                                 2015 July
## 4 Resort...
## 5 Resort...
                                 14
                                                 2015 July
    ... with 27 more variables: arrival date week number <dbl>,
       arrival date day of month <dbl>,
. . .
```

filter

filter to select a subset of rows

```
# bookings in City Hotels
hotels %>%
  filter(hotel == "City Hotel")
## # A tibble: 79,330 × 32
  hotel is_canceled lead_time arrival_date_ye... arrival_date_mo...
                  <dbl>
  <chr>
                            <dbl>
                                             <dbl> <chr>
## 1 City H...
                                              2015 July
## 2 City H...
                            88
                                              2015 July
## 3 City H...
                               65
                                              2015 July
## 4 City H...
                                              2015 July
                               92
## 5 City H...
                              100
                                              2015 July
## 6 City H...
                               79
                                              2015 July
## # ... with 79,324 more rows, and 27 more variables:
## # arrival_date_week_number <dbl>,
## # arrival_date_day_of_month <dbl>,
## # stays in weekend nights <dbl>, stays in week nights <dbl>,
      adults <dbl>, children <dbl>, babies <dbl>, meal <chr>,
## #
      country <chr>, market_segment <chr>,
## #
      distribution_channel <chr>, is_repeated_guest <dbl>, ...
## #
```



filter for many conditions at once

```
hotels %>%
  filter(
   adults == 0,
   children >= 1
   ) %>%
  select(adults, babies, children)
```

filter for more complex conditions

Logical operators in R

operator	definition	operator	definition
<	less than	x y	x OR y
<=	less than or equal to	is.na(x)	test if x is NA
>	greater than	!is.na(x)	test if x is not NA
>=	greater than or equal to	x %in% y	test if x is in y
==	exactly equal to	!(x %in% y)	test if x is not in y
!=	not equal to	!x	not x
x & y	x AND y		



distinct and count

distinct to filter for unique rows

... and arrange to order alphabetically

```
hotels %>%
    distinct(market_segment) %>%
    arrange(market_segment)

## # A tibble: 8 × 1

## market_segment

## <chr>
## 1 Aviation

## 2 Complementary

## 3 Corporate

## 4 Direct

## 5 Groups

## 6 Offline TA/TO

## 7 Online TA

## 8 Undefined
```

```
hotels %>%
  distinct(hotel, market_segment) %>%
  arrange(hotel, market_segment)

## # A tibble: 14 × 2

## hotel market_segment
```

```
<chr>
                  <chr>
   1 City Hotel
                  Aviation
   2 City Hotel
                  Complementary
   3 Citv Hotel
                  Corporate
   4 City Hotel
                  Direct
   5 City Hotel
                  Groups
   6 City Hotel
                  Offline TA/TO
   7 City Hotel
                  Online TA
                  Undefined
   8 City Hotel
   9 Resort Hotel Complementary
## 10 Resort Hotel Corporate
. . .
```

count to create frequency tables

```
# alphabetical order by default
hotels %>%
    count(market_segment)
```

```
## # A tibble: 8 × 2
##
    market segment
##
    <chr>
                    <int>
  1 Aviation
                      237
  2 Complementary
                     743
## 3 Corporate
                     5295
  4 Direct
                    12606
## 5 Groups
                    19811
## 6 Offline TA/TO 24219
## 7 Online TA
                    56477
## 8 Undefined
```

```
# descending frequency order
hotels %>%
   count(market_segment, sort = TRUE)
```

```
## # A tibble: 8 × 2
    market segment
    <chr>
                   <int>
## 1 Online TA
                   56477
## 2 Offline TA/T0 24219
## 3 Groups
                   19811
## 4 Direct
                   12606
## 5 Corporate
                    5295
## 6 Complementary
                     743
## 7 Aviation
                     237
## 8 Undefined
```

count and arrange

```
# ascending frequency order
hotels %>%
  count(market_segment) %>%
  arrange(n)
```

```
## # A tibble: 8 × 2
     market_segment
##
     <chr>
                    <int>
  1 Undefined
                      237
## 2 Aviation
                      743
  3 Complementary
                     5295
  4 Corporate
## 5 Direct
                    12606
## 6 Groups
                    19811
## 7 Offline TA/TO
                    24219
## 8 Online TA
                    56477
```

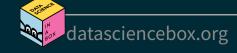
```
# descending frequency order
# just like adding sort = TRUE
hotels %>%
   count(market_segment) %>%
   arrange(desc(n))
```

```
## # A tibble: 8 × 2
     market_segment
    <chr>
                    <int>
## 1 Online TA
                    56477
## 2 Offline TA/TO
                    24219
## 3 Groups
                    19811
## 4 Direct
                    12606
## 5 Corporate
                     5295
                      743
## 6 Complementary
## 7 Aviation
                      237
## 8 Undefined
```

count for multiple variables

```
hotels %>%
  count(hotel, market_segment)
```

```
## # A tibble: 14 × 3
##
     hotel
                   market segment
                                      n
##
     <chr>
                   <chr>
                                  <int>
    1 City Hotel
                   Aviation
                                    237
                                    542
##
   2 City Hotel
                   Complementary
##
   3 City Hotel
                   Corporate
                                   2986
##
    4 City Hotel
                   Direct
                                   6093
   5 City Hotel
##
                   Groups
                                  13975
    6 City Hotel
                   Offline TA/TO
##
                                  16747
    7 City Hotel
                   Online TA
##
                                  38748
##
   8 City Hotel
                   Undefined
   9 Resort Hotel Complementary
                                    201
  10 Resort Hotel Corporate
                                   2309
  11 Resort Hotel Direct
                                   6513
  12 Resort Hotel Groups
                                   5836
  13 Resort Hotel Offline TA/TO
                                   7472
  14 Resort Hotel Online TA
                                  17729
```



order matters when you count

datasciencebox.org

```
# hotel type first
                                                 # market segment first
hotels %>%
                                                 hotels %>%
   count(hotel, market segment)
                                                   count(market segment, hotel)
## # A tibble: 14 × 3
                                                ## # A tibble: 14 × 3
##
      hotel
                   market segment
                                                ##
                                                      market segment hotel
                                      n
                                                                                       n
                                                      <chr>
##
      <chr>
                   <chr>
                                  <int>
                                                                     <chr>
                                                                                   <int>
                                                                                     237
##
    1 City Hotel
                   Aviation
                                    237
                                                    1 Aviation
                                                                     City Hotel
##
    2 City Hotel
                   Complementary
                                    542
                                                ##
                                                    2 Complementary
                                                                     City Hotel
                                                                                     542
##
    3 City Hotel
                   Corporate
                                   2986
                                                ##
                                                    3 Complementary
                                                                     Resort Hotel
                                                                                     201
    4 City Hotel
                   Direct
                                   6093
                                                    4 Corporate
                                                                     City Hotel
                                                                                    2986
##
                                                ##
    5 City Hotel
                                                    5 Corporate
                                                                                    2309
##
                   Groups
                                   13975
                                                ##
                                                                     Resort Hotel
                                                    6 Direct
                                                                                    6093
##
    6 City Hotel
                   Offline TA/TO
                                   16747
                                                ##
                                                                     City Hotel
    7 City Hotel
                   Online TA
                                                    7 Direct
                                                                     Resort Hotel
                                                                                    6513
##
                                   38748
##
   8 City Hotel
                   Undefined
                                                ##
                                                    8 Groups
                                                                     City Hotel
                                                                                   13975
                                                                                    5836
    9 Resort Hotel Complementary
                                    201
                                                ##
                                                    9 Groups
                                                                     Resort Hotel
  10 Resort Hotel Corporate
                                   2309
                                                   10 Offline TA/TO
                                                                     City Hotel
                                                                                   16747
  11 Resort Hotel Direct
                                   6513
                                                   11 Offline TA/TO
                                                                     Resort Hotel 7472
                                                   12 Online TA
                                                                     City Hotel
                                                                                   38748
  12 Resort Hotel Groups
                                   5836
  13 Resort Hotel Offline TA/TO
                                    7472
                                                   13 Online TA
                                                                     Resort Hotel 17729
  14 Resort Hotel Online TA
                                   17729
                                                   14 Undefined
                                                                      City Hotel
```

mutate

mutate to add a new variable

```
hotels %>%
  mutate(little_ones = children + babies) %>%
  select(children, babies, little_ones) %>%
  arrange(desc(little_ones))
```

Little ones in resort and city hotels

```
# Resort Hotel
hotels %>%
  mutate(little_ones = children + babies)
filter(
  little_ones >= 1,
  hotel == "Resort Hotel"
  ) %>%
  select(hotel, little_ones)
```

```
# City Hotel
hotels %>%
  mutate(little_ones = children + babies)
filter(
   little_ones >= 1,
   hotel == "City Hotel"
   ) %>%
select(hotel, little_ones)
```

What is happening in the following chunk?

```
hotels %>%
  mutate(little ones = children + babies) %>%
  count(hotel, little ones) %>%
  mutate(prop = n / sum(n))
## # A tibble: 12 × 4
                  little ones
##
     hotel
                                          prop
                        <dbl> <int>
                                         <dbl>
##
   <chr>
##
   1 City Hotel
                            0 73923 0.619
##
   2 City Hotel
                            1 3263 0.0273
##
   3 City Hotel
                              2056 0.0172
##
   4 City Hotel
                            3
                                 82 0.000687
   5 City Hotel
##
                                  1 0.00000838
##
   6 City Hotel
                           10
                                  1 0.00000838
##
   7 City Hotel
                                  4 0.0000335
##
   8 Resort Hotel
                            0 36131 0.303
##
   9 Resort Hotel
                               2183 0.0183
  10 Resort Hotel
                               1716 0.0144
  11 Resort Hotel
                            3 29 0.000243
  12 Resort Hotel
                           10
                                  1 0.00000838
```

datasciencebox.org

summarise and group_by

summarise for summary stats

```
# mean average daily rate for all bookings
hotels %>%
   summarise(mean_adr = mean(adr))
```

```
## # A tibble: 1 × 1
## mean_adr
## <dbl>
## 1 102.
```

summarise() changes the data frame entirely, it collapses rows down to a single summary statistic, and removes all columns that are irrelevant to the calculation.

summarise() also lets you get away with being sloppy and not naming your new column, but that's not recommended!



```
hotels %>%
  summarise(mean(adr))
```

```
## # A tibble: 1 × 1
## `mean(adr)`
## <dbl>
## 1 102.
```



```
hotels %>%
  summarise(mean_adr = mean(adr))
```

```
## # A tibble: 1 × 1
## mean_adr
## <dbl>
## 1 102.
```

group_by for grouped operations

```
# mean average daily rate for all booking at city and resort hotels
hotels %>%
   group_by(hotel) %>%
   summarise(mean_adr = mean(adr))
```

Calculating frequencies

The following two give the same result, so count is simply short for group_by then determine frequencies

```
hotels %>%
  group_by(hotel) %>%
  summarise(n = n())
```

```
hotels %>%
count(hotel)
```

Multiple summary statistics

summarise can be used for multiple summary statistics as well

```
hotels %>%
  summarise(
    min adr = min(adr),
    mean adr = mean(adr),
    median_adr = median(adr),
   max_adr = max(adr)
## # A tibble: 1 × 4
##
    min_adr mean_adr median_adr max_adr
    ##
## 1 -6.38 102. 94.6 5400
<!--
```

