

Data Transformation with dplyr

Data Analysis and Visualization (Fall 2019)

Instructor: Debopriya Ghosh

```
library(nycflights13)
```

```
library(tidyverse)
```

```
## — Attaching packages ————— tidyverse 1.2.1
```

```
## ✓ ggplot2 3.2.0      ✓ purrr  0.3.2
```

```
## ✓ tibble  2.1.3      ✓ dplyr  0.8.1
```

```
## ✓ tidyr   0.8.3      ✓ stringr 1.4.0
```

```
## ✓ readr   1.3.1      ✓ forcats 0.4.0
```

```
## — Conflicts ————— tidyverse_conflicts()
```

```
## ✗ dplyr::filter() masks stats::filter()
```

```
## ✗ dplyr::lag()     masks stats::lag()
```

```
data("flights")
```

```
flights
```

```
## # A tibble: 336,776 x 19
```

```
##   year month   day dep_time sched_dep_time dep_delay arr_time
```

```
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
```

```
## 1  2013     1     1     517             515           2     830
```

```
## 2  2013     1     1     533             529           4     850
```

```
## 3  2013     1     1     542             540           2     923
```

```
## 4  2013     1     1     544             545          -1    1004
```

```
## 5  2013     1     1     554             600          -6     812
```

```
## 6  2013     1     1     554             558          -4     740
```

```
## 7  2013     1     1     555             600          -5     913
```

```
## 8  2013     1     1     557             600          -3     709
```

```
## 9  2013     1     1     557             600          -3     838
```

```
## 10 2013     1     1     558             600          -2     753
```

```
## # ... with 336,766 more rows, and 12 more variables: sched_arr_time <int>,
```

```
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
```

```
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
```

```
## #   minute <dbl>, time_hour <dtm>
```

Filter rows with filter()

Subset observations based on their values.

```
Jan1 = filter(flights, month == 1, day == 1) # flights on January 1
```

```
(dec25 = filter(flights, month == 12, day == 25) )
```

```
## # A tibble: 719 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>   <int>
## 1  2013    12    25     456           500        -4     649
## 2  2013    12    25     524           515         9     805
## 3  2013    12    25     542           540         2     832
## 4  2013    12    25     546           550        -4    1022
## 5  2013    12    25     556           600        -4     730
## 6  2013    12    25     557           600        -3     743
## 7  2013    12    25     557           600        -3     818
## 8  2013    12    25     559           600        -1     855
## 9  2013    12    25     559           600        -1     849
## 10 2013    12    25     600           600         0     850
## # ... with 709 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>
```

Logical Operators

```
filter(flights, month == 11 | month == 12)
```

```
## # A tibble: 55,403 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>   <int>
## 1  2013    11     1         5          2359         6     352
## 2  2013    11     1        35          2250       105     123
## 3  2013    11     1       455           500        -5     641
## 4  2013    11     1       539           545        -6     856
## 5  2013    11     1       542           545        -3     831
## 6  2013    11     1       549           600       -11     912
## 7  2013    11     1       550           600       -10     705
## 8  2013    11     1       554           600        -6     659
## 9  2013    11     1       554           600        -6     826
## 10 2013    11     1       554           600        -6     749
## # ... with 55,393 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>
```

```
filter(flights, !(arr_delay >120 | dep_delay >120))
```

```
## # A tibble: 316,050 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>   <int>
## 1  2013     1     1     517           515         2     830
## 2  2013     1     1     533           529         4     850
## 3  2013     1     1     542           540         2     923
## 4  2013     1     1     544           545        -1    1004
## 5  2013     1     1     554           600        -6     812
## 6  2013     1     1     554           558        -4     740
```

```
## 7 2013 1 1 555 600 -5 913
## 8 2013 1 1 557 600 -3 709
## 9 2013 1 1 557 600 -3 838
## 10 2013 1 1 558 600 -2 753
## # ... with 316,040 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>
```

OR

```
filter(flights, arr_delay <= 120 , dep_delay <= 120)

## # A tibble: 316,050 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1 2013     1     1     517           515           2     830
## 2 2013     1     1     533           529           4     850
## 3 2013     1     1     542           540           2     923
## 4 2013     1     1     544           545          -1    1004
## 5 2013     1     1     554           600          -6     812
## 6 2013     1     1     554           558          -4     740
## 7 2013     1     1     555           600          -5     913
## 8 2013     1     1     557           600          -3     709
## 9 2013     1     1     557           600          -3     838
## 10 2013     1     1     558           600          -2     753
## # ... with 316,040 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>
```

Missing Values

```
df = tibble(x = c(1, NA, 3))
filter(df, x > 1) # excludes both FALSE and NA values
```

```
## # A tibble: 1 x 1
##       x
##   <dbl>
## 1     3
```

```
filter(df, is.na(x) | x > 1)
```

```
## # A tibble: 2 x 1
##       x
##   <dbl>
## 1    NA
## 2     3
```

EXERCISES

- (1) Find all flights that –
 - (i) Had an arrival delay of two or more hours

- (ii) Flew to Houston (IAH or HOU)
- (iii) Were operated by United, American, or Delta
- (iv) Departed in summer(July, August, September)
- (v) Arrived more than two hours late, but didn't leave late
- (vi) Departed between midnight and 6 am (inclusive)

(2) How many flights have a missing dep_time?

Arrange Rows with arrange()

`arrange(flights, year, month, day)`

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     1     517           515           2     830
## 2  2013     1     1     533           529           4     850
## 3  2013     1     1     542           540           2     923
## 4  2013     1     1     544           545          -1    1004
## 5  2013     1     1     554           600          -6     812
## 6  2013     1     1     554           558          -4     740
## 7  2013     1     1     555           600          -5     913
## 8  2013     1     1     557           600          -3     709
## 9  2013     1     1     557           600          -3     838
## 10 2013     1     1     558           600          -2     753
## # ... with 336,766 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>
```

`arrange(flights, desc(arr_delay))`

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     9     641           900        1301    1242
## 2  2013     6    15    1432          1935        1137    1607
## 3  2013     1    10    1121          1635        1126    1239
## 4  2013     9    20    1139          1845        1014    1457
## 5  2013     7    22     845          1600        1005    1044
## 6  2013     4    10    1100          1900         960    1342
## 7  2013     3    17    2321           810         911     135
## 8  2013     7    22    2257           759         898     121
## 9  2013    12     5     756          1700         896    1058
## 10 2013     5     3    1133          2055         878    1250
## # ... with 336,766 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,
```

```
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>

# missing values are sorted at the end
df = tibble(x = c(5,2,NA))
arrange(df,x)

## # A tibble: 3 x 1
##       x
##   <dbl>
## 1     2
## 2     5
## 3    NA

arrange(df,desc(x))

## # A tibble: 3 x 1
##       x
##   <dbl>
## 1     5
## 2     2
## 3    NA
```

EXERCISES

- (1) Sort flights to find the most delayed flights. Find flights that left earliest.
- (2) Sort flights to find the fastest flights.
- (3) Which flight traveled the longest? Which flight traveled the shortest?

Select Columns with select()

```
select(flights, year, month, day)
```

```
## # A tibble: 336,776 x 3
##   year month   day
##   <int> <int> <int>
## 1  2013     1     1
## 2  2013     1     1
## 3  2013     1     1
## 4  2013     1     1
## 5  2013     1     1
## 6  2013     1     1
## 7  2013     1     1
## 8  2013     1     1
## 9  2013     1     1
## 10 2013     1     1
## # ... with 336,766 more rows
```

```
select(flights, year:day)
```

```
## # A tibble: 336,776 x 3
##   year month   day
```

```
##      <int> <int> <int>
## 1  2013      1      1
## 2  2013      1      1
## 3  2013      1      1
## 4  2013      1      1
## 5  2013      1      1
## 6  2013      1      1
## 7  2013      1      1
## 8  2013      1      1
## 9  2013      1      1
## 10 2013      1      1
## # ... with 336,766 more rows
```

```
select(flights, -(year:day))
```

```
## # A tibble: 336,776 x 16
##   dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay
##   <int>         <int>         <dbl>   <int>         <int>         <dbl>
## 1     517           515           2     830           819           11
## 2     533           529           4     850           830           20
## 3     542           540           2     923           850           33
## 4     544           545          -1    1004          1022          -18
## 5     554           600          -6     812           837          -25
## 6     554           558          -4     740           728           12
## 7     555           600          -5     913           854           19
## 8     557           600          -3     709           723          -14
## 9     557           600          -3     838           846           -8
## 10    558           600          -2     753           745            8
## # ... with 336,766 more rows, and 10 more variables: carrier <chr>,
## #   flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
## #   distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
select(flights, time_hour, air_time, everything()) # moves specified columns
to beginning
```

```
## # A tibble: 336,776 x 19
##   time_hour          air_time year month   day dep_time sched_dep_time
##   <dtm>              <dbl> <int> <int> <int>   <int>         <int>
## 1 2013-01-01 05:00:00    227  2013     1     1     517           515
## 2 2013-01-01 05:00:00    227  2013     1     1     533           529
## 3 2013-01-01 05:00:00    160  2013     1     1     542           540
## 4 2013-01-01 05:00:00    183  2013     1     1     544           545
## 5 2013-01-01 06:00:00    116  2013     1     1     554           600
## 6 2013-01-01 05:00:00    150  2013     1     1     554           558
## 7 2013-01-01 06:00:00    158  2013     1     1     555           600
## 8 2013-01-01 06:00:00     53  2013     1     1     557           600
## 9 2013-01-01 06:00:00    140  2013     1     1     557           600
## 10 2013-01-01 06:00:00    138  2013     1     1     558           600
## # ... with 336,766 more rows, and 12 more variables: dep_delay <dbl>,
## #   arr_time <int>, sched_arr_time <int>, arr_delay <dbl>, carrier <chr>,
```

```
## # flight <int>, tailnum <chr>, origin <chr>, dest <chr>, distance <dbl>,
## # hour <dbl>, minute <dbl>

rename(flights, tail_num = tailnum) # rename variable

## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     1     517             515           2     830
## 2  2013     1     1     533             529           4     850
## 3  2013     1     1     542             540           2     923
## 4  2013     1     1     544             545          -1    1004
## 5  2013     1     1     554             600          -6     812
## 6  2013     1     1     554             558          -4     740
## 7  2013     1     1     555             600          -5     913
## 8  2013     1     1     557             600          -3     709
## 9  2013     1     1     557             600          -3     838
## 10 2013     1     1     558             600          -2     753
## # ... with 336,766 more rows, and 12 more variables: sched_arr_time <int>,
## #   arr_delay <dbl>, carrier <chr>, flight <int>, tail_num <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>
```

Add New Variables with mutate()

```
flights_sml = select(flights,
  year:day,
  ends_with("delay"),
  distance,
  air_time)

mutate(flights_sml,
  gain = arr_delay - dep_delay,
  speed = distance/air_time*60)

## # A tibble: 336,776 x 9
##   year month   day dep_delay arr_delay distance air_time gain speed
##   <int> <int> <int>     <dbl>     <dbl>     <dbl>   <dbl> <dbl> <dbl>
## 1  2013     1     1         2         11    1400     227     9   370.
## 2  2013     1     1         4         20    1416     227    16   374.
## 3  2013     1     1         2         33    1089     160    31   408.
## 4  2013     1     1        -1        -18    1576     183   -17   517.
## 5  2013     1     1        -6        -25     762     116   -19   394.
## 6  2013     1     1        -4         12     719     150    16   288.
## 7  2013     1     1        -5         19    1065     158    24   404.
## 8  2013     1     1        -3        -14     229      53   -11   259.
## 9  2013     1     1        -3         -8     944     140    -5   405.
## 10 2013     1     1        -2          8     733     138    10   319.
## # ... with 336,766 more rows
```

Now, we can refer to the column just created.

```
mutate(flights_sml,
       gain = arr_delay - dep_delay,
       hours = air_time / 60,
       gain_per_hour = gain/hours)

## # A tibble: 336,776 x 10
##   year month   day dep_delay arr_delay distance air_time  gain hours
##   <int> <int> <int>   <dbl>   <dbl>   <dbl>   <dbl> <dbl> <dbl>
## 1  2013     1     1         2        11    1400    227     9  3.78
## 2  2013     1     1         4        20    1416    227    16  3.78
## 3  2013     1     1         2        33    1089    160    31  2.67
## 4  2013     1     1        -1       -18    1576    183   -17  3.05
## 5  2013     1     1        -6       -25     762    116   -19  1.93
## 6  2013     1     1        -4        12     719    150    16  2.5
## 7  2013     1     1        -5        19    1065    158    24  2.63
## 8  2013     1     1        -3       -14     229     53   -11  0.883
## 9  2013     1     1        -3        -8     944    140    -5  2.33
## 10 2013     1     1        -2         8     733    138    10  2.3
## # ... with 336,766 more rows, and 1 more variable: gain_per_hour <dbl>
```

Keep only the new variables.

```
transmute(flights,
          gain = arr_delay - dep_delay,
          hours = air_time / 60,
          gain_per_hour = gain/hours )

## # A tibble: 336,776 x 3
##   gain hours gain_per_hour
##   <dbl> <dbl>   <dbl>
## 1     9  3.78         2.38
## 2    16  3.78         4.23
## 3    31  2.67        11.6
## 4   -17  3.05        -5.57
## 5   -19  1.93       -9.83
## 6    16  2.5         6.4
## 7    24  2.63         9.11
## 8   -11  0.883       -12.5
## 9     -5  2.33        -2.14
## 10   10  2.3         4.35
## # ... with 336,766 more rows
```

Grouped Summaries with summarize()

```
summarize(flights, delay = mean(dep_delay, na.rm = T))

## # A tibble: 1 x 1
##   delay
##   <dbl>
## 1  12.6
```



```

by_day = group_by(flights, year, month, day)
summarize(by_day, delay = mean(dep_delay, na.rm = T))

## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day delay
##   <int> <int> <int> <dbl>
## 1  2013     1     1  11.5
## 2  2013     1     2  13.9
## 3  2013     1     3  11.0
## 4  2013     1     4   8.95
## 5  2013     1     5   5.73
## 6  2013     1     6   7.15
## 7  2013     1     7   5.42
## 8  2013     1     8   2.55
## 9  2013     1     9   2.28
## 10 2013     1    10   2.84
## # ... with 355 more rows

```

Combining Multiple Operations with Pipe

1. Group flights by destination
2. Summarize to compute distance, avg. delay, and number of flights
3. Filter to remove noisy points and Honolulu airport, which is almost twice as far away to the next closest airport

```

delays = flights %>%
  group_by(dest) %>%
  summarize(
    count = n(),
    dist = mean(distance, na.rm = T),
    delay = mean(arr_delay, na.rm = T)
  ) %>%
  filter(count > 20, dest != "HNL")

```