

Data Transformation with dplyr

Ismail Fadeli

August 27, 2021

Data Transformation with dplyr

```
library(nycflights13)

library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --

## v ggplot2 3.3.5      v purrr 0.3.4
## v tibble 3.1.3       v dplyr 1.0.7
## v tidyr 1.1.3        v stringr 1.4.0
## v readr 2.0.1        v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

flights

## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517             515           2       830           819
## 2  2013     1     1     533             529           4       850           830
## 3  2013     1     1     542             540           2       923           850
## 4  2013     1     1     544             545          -1      1004          1022
## 5  2013     1     1     554             600          -6       812           837
## 6  2013     1     1     554             558          -4       740           728
## 7  2013     1     1     555             600          -5       913           854
## 8  2013     1     1     557             600          -3       709           723
## 9  2013     1     1     557             600          -3       838           846
## 10 2013     1     1     558             600          -2       753           745
## # ... with 336,766 more rows, and 11 more variables: 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>
```

If we want to see all the data in flights: View(flights)

Filter Rows with filter()

filter() allows you to subset observations based on their values. The first argument is the name of the dataframe. The second and the rest arguments are the expressions that filter the dataframe. For example, we can select all flights on January 1st with:

```
filter(flights, month == 1, day ==1)
```

```
## # A tibble: 842 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517             515           2     830           819
## 2  2013     1     1     533             529           4     850           830
## 3  2013     1     1     542             540           2     923           850
## 4  2013     1     1     544             545          -1    1004          1022
## 5  2013     1     1     554             600          -6     812           837
## 6  2013     1     1     554             558          -4     740           728
## 7  2013     1     1     555             600          -5     913           854
## 8  2013     1     1     557             600          -3     709           723
## 9  2013     1     1     557             600          -3     838           846
## 10 2013     1     1     558             600          -2     753           745
## # ... with 832 more rows, and 11 more variables: 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>
```

To assign the filtered result into a new dataframe we can use an assignment:

```
jan1 <- filter(flights, month == 1, day == 1)
```

If you want to print the results and save them in a variable at the same time, you can use parentheses:

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

```
## # A tibble: 719 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013    12    25     456             500          -4     649           651
## 2  2013    12    25     524             515           9     805           814
## 3  2013    12    25     542             540           2     832           850
## 4  2013    12    25     546             550          -4    1022          1027
## 5  2013    12    25     556             600          -4     730           745
## 6  2013    12    25     557             600          -3     743           752
## 7  2013    12    25     557             600          -3     818           831
## 8  2013    12    25     559             600          -1     855           856
## 9  2013    12    25     559             600          -1     849           855
## 10 2013    12    25     600             600           0     850           846
## # ... with 709 more rows, and 11 more variables: 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>
```

Comparisons

The standard operations that R uses are: >, >=, <, <=, != (not equal), and == (equal).

```
sqrt(2) ^ 2 == 2
```

```
## [1] FALSE
```

```
1/49 *49 == 1
```

```
## [1] FALSE
```

```
near(sqrt(2) ^ 2, 2)
```

```
## [1] TRUE
```

```
near(1 / 49 * 49, 1)
```

```
## [1] TRUE
```

Boolean operators in R are: & is “and”, | is “or”, ! is “not”. The following code shows the months of November or December:

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

```
## # A tibble: 55,403 x 19
```

```
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013    11     1       5           2359           6       352           345
## 2  2013    11     1      35           2250          105       123          2356
## 3  2013    11     1     455           500           -5       641           651
## 4  2013    11     1     539           545           -6       856           827
## 5  2013    11     1     542           545           -3       831           855
## 6  2013    11     1     549           600          -11       912           923
## 7  2013    11     1     550           600          -10       705           659
## 8  2013    11     1     554           600           -6       659           701
## 9  2013    11     1     554           600           -6       826           827
##10  2013    11     1     554           600           -6       749           751
```

```
## # ... with 55,393 more rows, and 11 more variables: 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>
```

```
nov_dec <- filter(flights, month %in% c(11, 12))
```

```
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 sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2       830           819
## 2  2013     1     1     533           529           4       850           830
## 3  2013     1     1     542           540           2       923           850
## 4  2013     1     1     544           545          -1      1004          1022
## 5  2013     1     1     554           600           -6       812           837
## 6  2013     1     1     554           558           -4       740           728
## 7  2013     1     1     555           600           -5       913           854
## 8  2013     1     1     557           600           -3       709           723
## 9  2013     1     1     557           600           -3       838           846
##10  2013     1     1     558           600           -2       753           745
```

```
## # ... with 316,040 more rows, and 11 more variables: 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 sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2       830           819
## 2  2013     1     1     533           529           4       850           830
## 3  2013     1     1     542           540           2       923           850
## 4  2013     1     1     544           545          -1      1004          1022
```

```
## 5 2013 1 1 554 600 -6 812 837
## 6 2013 1 1 554 558 -4 740 728
## 7 2013 1 1 555 600 -5 913 854
## 8 2013 1 1 557 600 -3 709 723
## 9 2013 1 1 557 600 -3 838 846
## 10 2013 1 1 558 600 -2 753 745
## # ... with 316,040 more rows, and 11 more variables: 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

if you want to determine if a value is missing, use `is.na()`:

```
x <- NA
```

```
is.na(x)
```

```
## [1] TRUE
```

```
df <- tibble(x = c(1, NA, 3))
filter(df, x > 1)
```

```
## # 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
```

Finding all flights that flew to Houston:

```
flights
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2     830           819
## 2  2013     1     1     533           529           4     850           830
## 3  2013     1     1     542           540           2     923           850
## 4  2013     1     1     544           545          -1    1004          1022
## 5  2013     1     1     554           600          -6     812           837
## 6  2013     1     1     554           558          -4     740           728
## 7  2013     1     1     555           600          -5     913           854
## 8  2013     1     1     557           600          -3     709           723
## 9  2013     1     1     557           600          -3     838           846
## 10 2013     1     1     558           600          -2     753           745
## # ... with 336,766 more rows, and 11 more variables: 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, dest == "IAH" | dest == "HOU")
```

```
## # A tibble: 9,313 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>       <dbl>   <int>         <int>
## 1  2013     1     1     517           515         2       830           819
## 2  2013     1     1     533           529         4       850           830
## 3  2013     1     1     623           627        -4       933           932
## 4  2013     1     1     728           732        -4      1041          1038
## 5  2013     1     1     739           739         0      1104          1038
## 6  2013     1     1     908           908         0      1228          1219
## 7  2013     1     1    1028          1026         2      1350          1339
## 8  2013     1     1    1044          1045        -1      1352          1351
## 9  2013     1     1    1114           900       134      1447          1222
## 10 2013     1     1    1205          1200         5      1503          1505
## # ... with 9,303 more rows, and 11 more variables: 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>
```

[View\(flights\)](#)

Finding all flights that were operated by United, American, or Delta:

```
filter(flights, carrier == "UA" | carrier == "AA" | carrier == "DL")
```

```
## # A tibble: 139,504 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>       <dbl>   <int>         <int>
## 1  2013     1     1     517           515         2       830           819
## 2  2013     1     1     533           529         4       850           830
## 3  2013     1     1     542           540         2       923           850
## 4  2013     1     1     554           600        -6       812           837
## 5  2013     1     1     554           558        -4       740           728
## 6  2013     1     1     558           600        -2       753           745
## 7  2013     1     1     558           600        -2       924           917
## 8  2013     1     1     558           600        -2       923           937
## 9  2013     1     1     559           600        -1       941           910
## 10 2013     1     1     559           600        -1       854           902
## # ... with 139,494 more rows, and 11 more variables: 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, between(x = month, left = 7, right = 9))
```

```
## # A tibble: 86,326 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>       <dbl>   <int>         <int>
## 1  2013     7     1         1          2029        212       236          2359
## 2  2013     7     1         2          2359         3       344           344
## 3  2013     7     1        29          2245        104       151           1
## 4  2013     7     1        43          2130        193       322           14
## 5  2013     7     1        44          2150        174       300           100
## 6  2013     7     1        46          2051        235       304          2358
## 7  2013     7     1        48          2001        287       308          2305
## 8  2013     7     1        58          2155        183       335           43
## 9  2013     7     1       100          2146        194       327           30
## 10 2013     7     1       100          2245        135       337          135
## # ... with 86,316 more rows, and 11 more variables: 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 Rows with Arrange()

Arrange() takes the dataframe and a set of column names to order by.

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

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>      <int>         <int>
## 1  2013     1     1     517           515         2        830           819
## 2  2013     1     1     533           529         4        850           830
## 3  2013     1     1     542           540         2        923           850
## 4  2013     1     1     544           545        -1       1004          1022
## 5  2013     1     1     554           600        -6        812           837
## 6  2013     1     1     554           558        -4        740           728
## 7  2013     1     1     555           600        -5        913           854
## 8  2013     1     1     557           600        -3        709           723
## 9  2013     1     1     557           600        -3        838           846
## 10 2013     1     1     558           600        -2        753           745
## # ... with 336,766 more rows, and 11 more variables: 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>
```

Use desc() to reorder by a column in descending order:

```
arrange(flights, desc(arr_delay))
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>      <int>         <int>
## 1  2013     1     9     641           900       1301       1242          1530
## 2  2013     6    15    1432          1935       1137       1607          2120
## 3  2013     1    10    1121          1635       1126       1239          1810
## 4  2013     9    20    1139          1845       1014       1457          2210
## 5  2013     7    22     845          1600       1005       1044          1815
## 6  2013     4    10    1100          1900        960       1342          2211
## 7  2013     3    17    2321           810        911        135          1020
## 8  2013     7    22    2257           759        898        121          1026
## 9  2013    12     5     756          1700        896       1058          2020
## 10 2013     5     3    1133          2055        878       1250          2215
## # ... with 336,766 more rows, and 11 more variables: 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 always 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
```

Select Columns with select()

```
# Select columns by name
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 all columns between year and day
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 all columns except those from year to day
select(flights, -(year:day))
```

```
## # A tibble: 336,776 x 16
##   dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier
##   <int>         <int>      <dbl>    <int>         <int>      <dbl> <chr>
## 1     517           515         2      830           819        11 UA
```

```
## 2      533      529      4      850      830      20 UA
## 3      542      540      2      923      850      33 AA
## 4      544      545     -1     1004     1022     -18 B6
## 5      554      600     -6      812      837     -25 DL
## 6      554      558     -4      740      728      12 UA
## 7      555      600     -5      913      854      19 B6
## 8      557      600     -3      709      723     -14 EV
## 9      557      600     -3      838      846      -8 B6
## 10     558      600     -2      753      745       8 AA
## # ... with 336,766 more rows, and 9 more variables: flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dtm>
```

There are a number of helper functions you can use within `select()`: `starts_with("abc")` matches names that begin with "abc". `ends_with("xyz")` matches names that end with "xyz". `contains("ijk")` matches names that contain "ijk". `matches("(.)\\1")` selects variables that match a regular expression. This one matches any variables that contain repeated characters. You'll learn more about regular expressions in Chapter 11. `num_range("x", 1:3)` matches `x1`, `x2`, and `x3`.

Use `rename()` to rename a variables:

```
rename(flights, tail_num = tailnum)
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2     830           819
## 2  2013     1     1     533           529           4     850           830
## 3  2013     1     1     542           540           2     923           850
## 4  2013     1     1     544           545          -1    1004          1022
## 5  2013     1     1     554           600          -6     812           837
## 6  2013     1     1     554           558          -4     740           728
## 7  2013     1     1     555           600          -5     913           854
## 8  2013     1     1     557           600          -3     709           723
## 9  2013     1     1     557           600          -3     838           846
## 10 2013     1     1     558           600          -2     753           745
## # ... with 336,766 more rows, and 11 more variables: 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>
```

Another option is to use `select()` in conjunction with the `everything()` helper. This is useful if we had a number of variables we wanted to move to the start of the dataframe.

```
select(flights, time_hour, air_time, everything())
```

```
## # 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>
```

Exercise:

1. Selecting dep_time, dep_delay, arr_time, arr_delay from flights:

```
select(flights, dep_time, dep_delay, arr_time, arr_delay)
```

```
## # A tibble: 336,776 x 4
##   dep_time dep_delay arr_time arr_delay
##   <int>     <dbl>   <int>     <dbl>
## 1     517         2     830         11
## 2     533         4     850         20
## 3     542         2     923         33
## 4     544        -1    1004        -18
## 5     554        -6     812        -25
## 6     554        -4     740         12
## 7     555        -5     913         19
## 8     557        -3     709        -14
## 9     557        -3     838         -8
## 10     558        -2     753          8
## # ... with 336,766 more rows
```

```
select(flights, dep_time, dep_delay, dep_delay)
```

```
## # A tibble: 336,776 x 2
##   dep_time dep_delay
##   <int>     <dbl>
## 1     517         2
## 2     533         4
## 3     542         2
## 4     544        -1
## 5     554        -6
## 6     554        -4
## 7     555        -5
## 8     557        -3
## 9     557        -3
## 10     558        -2
## # ... with 336,766 more rows
```

```
select(flights, contains("TIME"))
```

```
## # A tibble: 336,776 x 6
##   dep_time sched_dep_time arr_time sched_arr_time air_time time_hour
##   <int>         <int>   <int>         <int>     <dbl> <dtm>
## 1     517         515     830         819      227 2013-01-01 05:00:00
## 2     533         529     850         830      227 2013-01-01 05:00:00
## 3     542         540     923         850      160 2013-01-01 05:00:00
## 4     544         545    1004        1022      183 2013-01-01 05:00:00
## 5     554         600     812         837      116 2013-01-01 06:00:00
## 6     554         558     740         728      150 2013-01-01 05:00:00
## 7     555         600     913         854      158 2013-01-01 06:00:00
```

```
## 8      557      600      709      723      53 2013-01-01 06:00:00
## 9      557      600      838      846      140 2013-01-01 06:00:00
## 10     558      600      753      745      138 2013-01-01 06:00:00
## # ... with 336,766 more rows
```

Adding new variables with `mutate()`