

Lecture Assignment 6

Taiki Yamashita

2024-04-23

5.2.4 Exercises

```
# importing libraries...
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.5.0      v tibble    3.2.1
## v lubridate  1.9.3      v tidyr     1.3.1
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors

library(nycflights13)
```

1) Find all flights that...

Had an arrival delay of two or more hours

```
# since the arr_delay variable is measured in minutes, we will find flights with an arrival delay of 120
filter(flights, arr_delay >= 120)
```

```
## # A tibble: 10,200 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       811             630          101      1047             830
## 2  2013     1     1       848             1835         853      1001             1950
## 3  2013     1     1       957             733          144      1056             853
## 4  2013     1     1      1114             900          134      1447             1222
## 5  2013     1     1      1505             1310         115      1638             1431
## 6  2013     1     1      1525             1340          105      1831             1626
## 7  2013     1     1      1549             1445           64      1912             1656
## 8  2013     1     1      1558             1359          119      1718             1515
## 9  2013     1     1      1732             1630           62      2028             1825
```

```
## 10 2013      1      1      1803          1620          103      2008          1750
## # i 10,190 more rows
## # i 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>
```

Flew to Houston(IAH or HOU)

```
# the flights that flew to Houston are those flights where the destination (dest) is either "IAH" or "HOU"
filter(flights, dest %in% c("IAH", "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
## # i 9,303 more rows
## # i 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>
```

Were operated by United, American, or Delta

```
# using %in% here would make it more compact
filter(flights, carrier %in% c("AA", "DL", "UA"))
```

```
## # 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
## # i 139,494 more rows
## # i 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>
```

Departed in summer (July, August, September)

```
# the %in% operator is an alternative. we can use : here to specify the integer range.
filter(flights, month %in% 7:9)
```

```
## # 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
## # i 86,316 more rows
## # i 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>
```

Arrived more than two hours late, but didn't leave late

```
# flights that arrived more than 120 minutes late, but didn't leave late, dep_delay which represents de
filter(flights, arr_delay > 120, dep_delay <= 0)
```

```
## # A tibble: 29 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    27    1419           1420          -1     1754           1550
## 2  2013    10     7    1350           1350           0     1736           1526
## 3  2013    10     7    1357           1359          -2     1858           1654
## 4  2013    10    16     657            700          -3     1258           1056
## 5  2013    11     1     658            700          -2     1329           1015
## 6  2013     3    18    1844           1847          -3         39           2219
## 7  2013     4    17    1635           1640          -5     2049           1845
## 8  2013     4    18     558            600          -2     1149            850
## 9  2013     4    18     655            700          -5     1213            950
## 10 2013     5    22    1827           1830          -3     2217           2010
## # i 19 more rows
## # i 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>
```

Were delayed by at least an hour, but made up over 30 minutes of flight

```
filter(flights, dep_delay >= 60, (dep_delay - arr_delay > 30))
```

```
## # A tibble: 1,844 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    2205           1720        285     46           2040
## 2  2013     1     1    2326           2130        116    131           18
## 3  2013     1     3    1503           1221        162   1803          1555
## 4  2013     1     3    1839           1700         99   2056          1950
## 5  2013     1     3    1850           1745         65   2148          2120
## 6  2013     1     3    1941           1759        102   2246          2139
## 7  2013     1     3    1950           1845         65   2228          2227
## 8  2013     1     3    2015           1915         60   2135          2111
## 9  2013     1     3    2257           2000        177     45          2224
## 10 2013     1     4    1917           1700        137   2135          1950
## # i 1,834 more rows
## # i 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>
```

Departed between midnight and 6am (inclusive)

```
filter(flights, dep_time >= 2400 | dep_time <= 600)
```

```
## # A tibble: 9,373 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
## # i 9,363 more rows
## # i 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>
```

2) Another useful dplyr filtering helper is `between()`. What does it do? Can you use it to simplify the code needed to answer the previous challenges.

```
# the expression between(x, left, right) is the same as saying that x >= left & x <= right.
# therefore, for the question that flights that departed during the summer, it can be rewritten in this
filter(flights, between(month, 7, 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
## # i 86,316 more rows
## # i 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>
```

3) How many flights have a missing dep_time? What other variables are missing? What might these rows represent?

```
filter(flights, is.na(dep_time)) #using the is.na() function will generate me rows of flights with a mi.
```

```
## # A tibble: 8,255 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      NA           1630        NA      NA           1815
## 2  2013     1     1      NA           1935        NA      NA           2240
## 3  2013     1     1      NA           1500        NA      NA           1825
## 4  2013     1     1      NA           600         NA      NA           901
## 5  2013     1     2      NA           1540        NA      NA           1747
## 6  2013     1     2      NA           1620        NA      NA           1746
## 7  2013     1     2      NA           1355        NA      NA           1459
## 8  2013     1     2      NA           1420        NA      NA           1644
## 9  2013     1     2      NA           1321        NA      NA           1536
## 10 2013     1     2      NA           1545        NA      NA           1910
## # i 8,245 more rows
## # i 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>
```

```
summary(flights) # the summary() function has the number of missing values for all of the non-character
```

```
##      year      month      day      dep_time      sched_dep_time
## Min.   :2013   Min.    : 1.000   Min.    : 1.00   Min.     : 1    Min.     : 106
## 1st Qu.:2013   1st Qu.: 4.000   1st Qu.: 8.00   1st Qu.: 907    1st Qu.: 906
## Median :2013   Median : 7.000   Median :16.00   Median :1401    Median :1359
## Mean   :2013   Mean    : 6.549   Mean     :15.71   Mean     :1349    Mean     :1344
## 3rd Qu.:2013   3rd Qu.:10.000   3rd Qu.:23.00   3rd Qu.:1744    3rd Qu.:1729
## Max.    :2013   Max.     :12.000   Max.     :31.00   Max.     :2400    Max.     :2359
##                                     NA's      :8255
```

```
##      dep_delay      arr_time      sched_arr_time      arr_delay
## Min.   : -43.00   Min.    :    1   Min.    :    1   Min.    : -86.000
## 1st Qu.:  -5.00   1st Qu.:1104   1st Qu.:1124   1st Qu.: -17.000
## Median :  -2.00   Median :1535   Median :1556   Median :  -5.000
## Mean   :  12.64   Mean    :1502   Mean    :1536   Mean    :   6.895
## 3rd Qu.:  11.00   3rd Qu.:1940   3rd Qu.:1945   3rd Qu.:  14.000
## Max.   :1301.00   Max.    :2400   Max.    :2359   Max.    :1272.000
## NA's   :8255     NA's    :8713                     NA's    :9430
##      carrier      flight      tailnum      origin
## Length:336776   Min.    :    1   Length:336776   Length:336776
## Class :character 1st Qu.: 553   Class :character Class :character
## Mode  :character Median :1496   Mode  :character Mode  :character
##                      Mean    :1972
##                      3rd Qu.:3465
##                      Max.    :8500
##
##      dest      air_time      distance      hour
## Length:336776   Min.    : 20.0   Min.    : 17   Min.    : 1.00
## Class :character 1st Qu.: 82.0   1st Qu.: 502   1st Qu.: 9.00
## Mode  :character Median :129.0   Median : 872   Median :13.00
##                      Mean    :150.7   Mean    :1040   Mean    :13.18
##                      3rd Qu.:192.0   3rd Qu.:1389   3rd Qu.:17.00
##                      Max.    :695.0   Max.    :4983   Max.    :23.00
##                      NA's    :9430
##      minute      time_hour
## Min.    : 0.00   Min.    :2013-01-01 05:00:00.00
## 1st Qu.: 8.00   1st Qu.:2013-04-04 13:00:00.00
## Median :29.00   Median :2013-07-03 10:00:00.00
## Mean    :26.23   Mean    :2013-07-03 05:22:54.64
## 3rd Qu.:44.00   3rd Qu.:2013-10-01 07:00:00.00
## Max.    :59.00   Max.    :2013-12-31 23:00:00.00
##
```

5.3.1 Exercises

1) How could you use `arrange()` to sort all missing values to the start? (Hint: use `is.na()`).

```
# what the arrange() functions does is it puts NA values last.
# instead to put th NA first, we can add something that checks if the column is missing a value.
arrange(flights, desc(is.na(dep_time)), dep_time)
```

```
## # 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     NA             1630             NA         NA             1815
## 2  2013     1     1     NA             1935             NA         NA             2240
## 3  2013     1     1     NA             1500             NA         NA             1825
## 4  2013     1     1     NA              600             NA         NA              901
## 5  2013     1     2     NA             1540             NA         NA             1747
## 6  2013     1     2     NA             1620             NA         NA             1746
```

```
## 7 2013 1 2 NA 1355 NA NA 1459
## 8 2013 1 2 NA 1420 NA NA 1644
## 9 2013 1 2 NA 1321 NA NA 1536
## 10 2013 1 2 NA 1545 NA NA 1910
## # i 336,766 more rows
## # i 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>
```

since desc(is.na(dep_time)) is TRUE when dep_time is missing and FALSE when it's not, the rows with th

What we are basically saying is that those which are 'TRUE' to being 'NA', sort them in descending order.

5.4.1 Exercises

1) Brainstorm as many ways as possible to select dep_time, dep_delay, arr_time, and arr_delay from flights.

```
vars <- c("dep_time", "dep_delay", "arr_time", "arr_delay")
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
## # i 336,766 more rows
```

```
select(flights, starts_with("dep"), starts_with("arr"))
```

```
## # 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
## # i 336,766 more rows
```

```
select(flights, one_of(vars))
```

```
## # 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
## # i 336,766 more rows
```

```
select_(flights, .dots = vars)
```

```
## Warning: 'select_()' was deprecated in dplyr 0.7.0.
## i Please use 'select()' instead.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```

```
## # 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
## # i 336,766 more rows
```

```
select_(flights, "dep_time", "dep_delay", "arr_time", "arr_delay")
```

```
## Warning: 'select_()' was deprecated in dplyr 0.7.0.
## i Please use 'select()' instead.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```

```
## # 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
## # i 336,766 more rows
```

```
select(flights, matches("dep"), matches("arr"), -matches("sched"), -carrier)
```

```
## # 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
## # i 336,766 more rows
```

```
select(flights, contains("dep"), contains("arr"), -contains("sched"), -carrier)
```

```
## # 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
## # i 336,766 more rows
```

```
select(flights, matches("^dep|^arr"))
```

```
## # 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
## # i 336,766 more rows
```

```
select(flights, matches("time$|delay$"), -contains("sched"), -contains("air"))
```

```
## # 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
## # i 336,766 more rows
```

```
select(flights, matches("^dep|arr_delay|time$"))
```

```
## # A tibble: 336,776 x 7
##   dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay air_time
##   <int>         <int>     <dbl>   <int>         <int>     <dbl>   <dbl>
## 1      517         515         2      830           819        11      227
## 2      533         529         4      850           830        20      227
## 3      542         540         2      923           850        33      160
## 4      544         545        -1     1004          1022       -18      183
## 5      554         600        -6      812           837       -25      116
## 6      554         558        -4      740           728        12      150
## 7      555         600        -5      913           854        19      158
## 8      557         600        -3      709           723       -14        53
## 9      557         600        -3      838           846        -8      140
## 10     558         600        -2      753           745         8      138
## # i 336,766 more rows
```