Richard Gresham

CPSC 375 Big Data

September 27, 2022

Homework #4

Prepare your answers as a single PDF file.

Group work: You may work in groups of 1-3. Include all group member names in the PDF file.

Only one person in the group should submit to Canvas.

Due: check on Canvas.

Load the nycflights13 library (will have to install the nycflights13 package first) which contains

flight arrival and departure data in a table called flights. Apply the tidyverse's data wrangling

verbs to answer these questions. For each question, give only the (one line as a data pipeline)

code beginning with flights %>%

1. List data only for flights that departed on February 12, 2013. Input: flights %>%

filter(year == 2013, month == 2, day == 12)

```
# A tibble: 893 × 19
                  day dep_time sched_d...1
    year month
   <int> <int> <int>
                                      <int>
                          <int>
    2013
              2
                   12
                              17
                                      2245
              2
 2
    2013
                    12
                             506
                                        500
 3
    2013
              2
                   12
                             520
                                        525
 4
              2
                   12
                                        530
    2013
                             524
 5
              2
    2013
                   12
                             535
                                        540
 6
              2
    2013
                   12
                             539
                                        540
 7
              2
    <u>2</u>013
                   12
                            551
                                        600
 8
              2
                   12
    2013
                             552
                                        600
              2
 9
    2013
                   12
                             553
                                        600
10
              2
                   12
                            555
    2013
                                        600
# ... with 883 more rows, 14 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>, ...
# i Use `print(n = ...)` to see more rows, and `colnames()` to s
ee all variable names
```

2. List data only for flights that were delayed (both arrival and departure) by more than 2

hours.

Input: flights %>% filter(dep_delay > 120, arr_delay > 120)

```
> flights %>% filter(dep_delay > 120, arr_delay >120)
# A tibble: 8.335 x 19
```

" A CIBBIC: 0,555 A 15									
	year	month	day	dep_time	sched¹	dep_d²	arr_t…³	sched…⁴	arr_d…⁵
	<int></int>	<int></int>	<int></int>	<int></int>	<int></int>	<db7></db7>	<int></int>	<int></int>	<db7></db7>
1	<u>2</u> 013	1	1	848	<u>1</u> 835	853	<u>1</u> 001	<u>1</u> 950	851
2	<u>2</u> 013	1	1	957	733	144	<u>1</u> 056	853	123
3	<u>2</u> 013	1	1	<u>1</u> 114	900	134	<u>1</u> 447	<u>1</u> 222	145
4	<u>2</u> 013	1	1	<u>1</u> 815	<u>1</u> 325	290	<u>2</u> 120	<u>1</u> 542	338
5	<u>2</u> 013	1	1	<u>1</u> 842	<u>1</u> 422	260	<u>1</u> 958	<u>1</u> 535	263
6	<u>2</u> 013	1	1	<u>1</u> 856	<u>1</u> 645	131	<u>2</u> 212	<u>2</u> 005	127
7	<u>2</u> 013	1	1	<u>1</u> 934	<u>1</u> 725	129	<u>2</u> 126	<u>1</u> 855	151
8	2013	1	1	<u>1</u> 938	<u>1</u> 703	155	2109	<u>1</u> 823	166
9	<u>2</u> 013	1	1	<u>1</u> 942	<u>1</u> 705	157	<u>2</u> 124	<u>1</u> 830	174
10	<u>2</u> 013	1	1	<u>2</u> 006	<u>1</u> 630	216	<u>2</u> 230	<u>1</u> 848	222
44		0 225	mono I	10 ×	vono Voni	ablact co	- nnion	ch n>	

- # ... with 8,325 more rows, 10 more variables: carrier <chr>,
- # flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
 # air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
- # time_hour <dttm>, and abbreviated variable names 'sched_dep_time,
- # ²dep_delay, ³arr_time, ⁴sched_arr_time, ⁵arr_delay
- # i Use `print(n = ...)` to see more rows, and `colnames()` to see all variable names

3. List data only for flights that were delayed (either arrival or departure) by more than 2 hours.

Input: flights %>% filter(dep_delay > 120 | arr_delay > 120)

Output:

```
# A tibble: 11,422 × 19
                      day dep_time sched_...¹ dep_d...² arr_t...³ sched...⁴ arr_d...⁵
     year month
                                                         <db7>
     <int> <int> <int>
                                 <int>
                                              <int>
                                                                    <int>
                                                                               <int>
     <u>2</u>013
                 1
                          1
                                   811
                                               630
                                                           101
                                                                     <u>1</u>047
                                                                                  830
                                                                                             137
     <u>2</u>013
                 1
                                    848
                                               <u>1</u>835
                                                           853
                                                                     <u>1</u>001
                                                                                <u>1</u>950
                                                                                             851
     <del>2</del>013
                                    957
                                                733
                                                           144
                                                                     1056
                                                                                  853
                                                                                             123
                 1
                          1
                                                                     <u>1</u>447
     <u>2</u>013
                                                900
                                                            134
                 1
                          1
                                  <u>1</u>114
                                                                                <u>1</u>222
                                                                                             145
     <u>2</u>013
                 1
                          1
                                   <u>1</u>505
                                               <u>1</u>310
                                                            115
                                                                      <u>1</u>638
                                                                                 <u>1</u>431
                                                                                             127
     <u>2</u>013
                                  <u>1</u>525
                                               <u>1</u>340
                                                           105
                                                                      <u>1</u>831
                                                                                <u>1</u>626
                                                                                             125
                                               <u>1</u>338
     <u>2</u>013
                 1
                          1
                                  <u>1</u>540
                                                            122
                                                                      <u>2</u>020
                                                                                 1825
                                                                                             115
     <u>2</u>013
                          1
                                   <u>1</u>549
                                               <u>1</u>445
                                                            64
                                                                      <u>1</u>912
                                                                                 <u>1</u>656
                 1
                                                                                             136
                                               <u>1</u>359
                                                            119
     <u>2</u>013
                 1
                          1
                                   <u>1</u>558
                                                                      <u>1</u>718
                                                                                 <u>1</u>515
                                                                                             123
     <u>2</u>013
                          1
                                   <u>1</u>732
                                               <u>1</u>630
                                                             62
                                                                      <u>2</u>028
                                                                                 <u>1</u>825
                                                                                             123
# ... with 11,412 more rows, 10 more variables: carrier <chr>,
    flight <int>, tailnum <chr>, origin <chr>, dest <chr>
     air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
    time_hour <dttm>, and abbreviated variable names 'sched_dep_time,
     ^{2}dep_delay, ^{3}arr_time, ^{4}sched_arr_time, ^{5}arr_delay
\# i Use `print(n = ...)` to see more rows, and `colnames()` to see all va
riable names
```

As you can see above at times of early departures we see the data still because the arr_delay was greater than 2 hours.

4. List data only for flights that were operated by United, American, or Delta.

These need to be looked via their carrier codes United = "UA", American = "AA" Delta = "DL"

Input: flights %>% filter(carrier == "UA" | carrier == "AA" | carrier == "DL")

```
# A tibble: 139,504 x 19
    year month
                  day dep_time sched_d...1
    <int> <int> <int>
                          <int>
                    1
                            517
                                        515
    2013
              1
    2013
              1
                     1
                            533
                                        529
 3
    2013
              1
                     1
                            542
                                        540
 4
    2013
                            554
                                        600
              1
                     1
 5
    2013
              1
                     1
                            554
                                        558
 6
    2013
              1
                     1
                            558
                                        600
    2013
              1
                     1
                            558
                                        600
 8
                     1
                                        600
    2013
              1
                            558
    2013
 9
              1
                     1
                            559
                                        600
10
    2013
              1
                     1
                            559
                                        600
 ... with 139,494 more rows, 14 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>, ...
# i Use \hat{} print(n = ...) to see more rows, and \hat{} colnames() to see all variable names
```

5. Sort data in order of fastest flights (air_time).

Arrange function automatically does it in ascending order.

Input: flights %>% arrange(-air_time)

```
# A tibble: 336.776 x 19
                   day dep_time sched_...¹ dep_d...² arr_t...³ sched...⁴ arr_d...⁵
    year month
   <int> <int>
                  <int>
                            <int>
                                       <int>
                                                <db7>
                                                          <int>
                                                                   <int>
                                                                            <db7>
               3
                     17
                             1337
                                        1335
                                                                    1836
    2013
                                                     2
                                                          1937
                                                                                61
               2
    2013
                      6
                              853
                                         900
                                                   -7
                                                          1542
                                                                    1540
                                                                                 2
                                        <u>1</u>000
    2013
               3
                                                                                21
                     15
                             1001
                                                    1
                                                           1551
                                                                    1530
               3
 4
    2013
                     17
                             1006
                                        1000
                                                          1607
                                                                    1530
                                                                                37
                                                     6
               3
 5
    2013
                     16
                             <u>1</u>001
                                        1000
                                                     1
                                                          1544
                                                                    1530
                                                                                14
               2
 6
    2013
                      5
                              900
                                         900
                                                    0
                                                          1555
                                                                    1540
                                                                                15
 7
                                         930
                                                     6
    2013
              11
                     12
                              936
                                                          1630
                                                                    1530
                                                                                60
 8
    2013
               3
                     14
                              958
                                        1000
                                                   -2
                                                          1542
                                                                                12
                                                                    1530
 9
    2013
                     20
                             1006
                                                     6
                                                                                44
              11
                                        1000
                                                           1639
                                                                    1555
10
    2013
               3
                     15
                             1342
                                       1335
                                                     7
                                                          1924
                                                                    1836
                                                                                48
 ... with 336,766 more rows, 10 more variables: carrier <chr>,
    flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
    air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
    time_hour <dttm>, and abbreviated variable names 'sched_dep_time,
    <sup>2</sup>dep_delay, <sup>3</sup>arr_time, <sup>4</sup>sched_arr_time, <sup>5</sup>arr_delay
# i Use `print(n = ...)` to see more rows, and `colnames()` to see all va
riable names
```

6. Sort data in order of longest duration flights (air_time).

Longest to shortest is in descending order.

Input: flights %>% arrange(desc(-air_time))

Output:

```
# A tibble: 336,776 \times 19
                    day dep_time sched_...¹ dep_d...² arr_t...³ sched...⁴ arr_d...⁵
    year month
    <int> <int> <int>
                            <int>
                                                               <int>
                                       <int>
                                                     <db7>
                                                                          <int>
                                                                                    <db7>
    2013
              1 16
                                <u>1</u>355
                                           <u>1</u>315
                                                        40
                                                                1442
                                                                           1411
                                                                                        31
     <u>2</u>013
               4
                       13
                                 537
                                             527
                                                        10
                                                                 622
                                                                            628
                                                                                        -6
                     6
    <u>2</u>013
               12
                                 922
                                            851
                                                        31
                                                                1021
                                                                            954
                                                                                        27
    <u>2</u>013
                      3
                                                        24
               2
                                <u>2</u>153
                                            <u>2</u>129
                                                                <u>2</u>247
                                                                           2224
                                                                                        23
     <u>2</u>013
                2
                       5
                                <u>1</u>303
                                            <u>1</u>315
                                                       -12
                                                                <u>1</u>342
                                                                           <u>1</u>411
                                                                                       -29
 6
     <u>2</u>013
                2
                       12
                                <u>2</u>123
                                            <u>2</u>130
                                                        -7
                                                                           <u>2</u>225
                                                                                       -14
                                                                <u>2</u>211
     <u>2</u>013
                3
                      2
                                <u>1</u>450
                                            <u>1</u>500
                                                       -10
                                                                <u>1</u>547
                                                                           1608
                                                                                       -21
 8
     2013
                3
                       8
                                <u>2</u>026
                                            <u>1</u>935
                                                         51
                                                                 <u>2</u>131
                                                                           <u>2</u>056
                                                                                        35
    <u>2</u>013
                                <u>1</u>456
                                            <u>1</u>329
                                                                           <u>1</u>426
                3
                       18
                                                         87
                                                                 <u>1</u>533
                                                                                        67
                                <u>2</u>226
                                            <u>2</u>145
10
    2013
                3
                       19
                                                         41
                                                                 <u>2</u>305
                                                                           <u>2</u>246
                                                                                        19
# ... with 336,766 more rows, 10 more variables: carrier <chr>,
     flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
     air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
     time_hour <dttm>, and abbreviated variable names ¹sched_dep_time,
    2dep_delay, 3arr_time, 4sched_arr_time, 5arr_delay
# i Use `print(n = ...)` to see more rows, and `colnames()` to see all va
riable names
```

7. Show only the origin and destination of flights sorted by longest flights.

In this one arrange must be used first as select will limit our observations to only columns dest or origin.

Input: flights %>% arrange(desc(-air_time)) %>% select(origin, dest)

```
# A tibble: 336,776 \times 2
   origin dest
   <chr>
           <chr>
 1 \text{ EWR}
           BDL
 2 EWR
           BDL
 3 EWR
           BDL
 4 EWR
           PHL
 5 EWR
           BDL
 6 EWR
           PHL
 7 LGA
           BOS
 8 JFK
           PHL
 9 EWR
           BDL
10 EWR
           BDL
# ... with 336,766 more rows
# i Use `print(n = ...) ` to see more rows
```

8. Add a new variable that indicates the total delay (both departure and arrival delay).

```
Input: flights %>% mutate(totaldelay = dep_delay + arr_delay)
```

Output: to show any real output I will use an example where I output totaldelay,

dep_delay, and arr_delay.

flights %>% mutate(totaldelay = dep_delay + arr_delay) %>% select(dep_delay, arr_delay, totaldelay)

To permanently keep our new total delay we must save the output back into flights.

9. Show only the origin and destination of flights sorted by descending order of total delay.

```
Input: flights %>% mutate(totaldelay = dep_delay + arr_delay) %>% arrange(desc(totaldelay)) %>% select(origin, dest)
```

```
# A tibble: 336,776 \times 2
   origin dest
   <chr>
            <chr>
 1 \text{ EWR}
            SF0
 2 JFK
            SF<sub>0</sub>
 3 LGA
            MSY
 4 JFK
            PHX
 5 JFK
            MKE
 6 EWR
            SF<sub>0</sub>
 7 EWR
            SEA ·
 8 JFK
            MCI
 9 LGA
            DEN
10 LGA
            DSM
# ... with 336,766 more rows
# i Use `print(n = ...)` to see more rows
```

10. Show only the origin and destination of 10 most delayed flights [Hint: there are multiple ways of solving this. Some additional functions that you will find useful are head(), slice(), min_rank().]

Input: flights %>% mutate(totaldelay = dep_delay + arr_delay) %>% arrange(desc(totaldelay)) %>% select(origin, dest) %>% slice(1:10) Output:

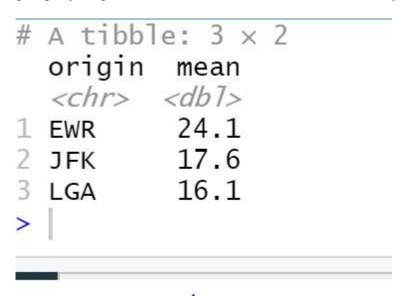
```
# A tibble: 10 \times 2
   origin dest
    <chr>
            <chr>
 1 \text{ EWR}
            SF0
 2 JFK
            SF0
 3 LGA
            MSY
 4 JFK
            PHX
 5 JFK
            MKE
 6 EWR
            SF<sub>0</sub>
 7 EWR
            SEA
 8 JFK
            MCI
 9 LGA
            DEN
10 LGA
            DSM
```

11. Show the average total delay for all flights.

Input: flights %>% mutate(totaldelay = dep_delay + arr_delay) %>% summarize(mean = mean(totaldelay, na.rm = TRUE)) Output:

12. Show the average total delay for every departure city.

Input: flights %>% mutate(totaldelay = dep_delay + arr_delay) %>% group_by(origin) %>% summarize(mean = mean(totaldelay, na.rm = TRUE)) Output:



13. Show the average total delay for every departure-arrival city pair.

Input: flights %>% mutate(totaldelay = dep_delay + arr_delay) %>% group_by(origin, dest) %>% summarize(mean = mean(totaldelay, na.rm = TRUE))