# Exercise 2: Reporting, Data Wrangling and Graphing

### Jianhui Gao

#### 07/08/2025

- Quick R
- Rstudio cheatsheet
- Rstudio for beginners

## Part 1: Analyze NYC flight delays.

Install the "nycflits13" package. The data comes from the US Bureau of Transportation Statistics. Using the data, complete the following tasks:

- 1. Find all flights that had an arrival delay of >4 hours, return the first 5 row. (Note: arr\_delay is in mins)
- 2. Find all flight names that flew from JFK to IAH, i.e. return only unique values of "flight" variable after filtering. Hint: unique() would help.
- 3. Find how many flights were operated by UA.
- 4. Find how many unique flights were operated by UA.
- 5. Sort flights that have the most delayed flights. Show the first 5 row.
- 6. Generate a scatter plot with x-axis dist and y-axis delay, where each dot is a unique flights and destination, dist is the average distance of each destination dest, and delay is the average delay time arr\_delay, with the size of dot equals to the count of delay records.

819

830

850

837

728

1022

```
library(tidyverse)
library(nycflights13)
head(flights)
```

```
## # A tibble: 6 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
## 2
      2013
                1
                              533
                                               529
                                                            4
                                                                    850
                       1
                                                            2
## 3
      2013
                              542
                                               540
                                                                    923
                1
                       1
## 4
      2013
                              544
                                               545
                                                                   1004
                1
                                                           -1
                       1
## 5
      2013
                1
                       1
                              554
                                               600
                                                           -6
                                                                    812
## 6
      2013
                1
                       1
                              554
                                               558
                                                           -4
                                                                    740
## # 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 <dttm> ## #

#### Solution

1. Find all flights that had an arrival delay of >4 hours, i.e. return the first 5 row. (Note: arr\_delay is in mins)

flights %>% filter(arr\_delay > 240) %>% head(5)

## # A tibble: 5 x 19

## day dep\_time sched\_dep\_time dep\_delay arr\_time sched\_arr\_time year month ## <int> <int> <int> <int> <int> <dbl> <int> <int> 2013 ## 1 1 1 848 1835 853 1001 1950 ## 2 2013 1815 1325 290 2120 1542 1 1 ## 3 2013 1 1 1842 1422 260 1958 1535 ## 4 2013 1700 255 1 1 2115 2330 1920 ## 5 2013 1 1 2205 1720 285 46 2040 ## # 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 <dttm>

2. Find all flight names that flew from JFK to IAH, i.e. return only unique values of "flight" variable after filtering. Hint: unique() would help.

df <- flights %>% filter(origin == "JFK" & dest == "IAH")
unique(df\$flight)

## [1] 211 1901 523

3. Find how many flights were operated by UA.

```
nrow(filter(flights, carrier %in% c("UA")))
```

## [1] 58665

4. Find how many unique flights were operated by UA.

```
df <- filter(flights, carrier %in% c("UA"))
length(unique(df$flight))</pre>
```

## [1] 1285

5. Sort flights that have the most delayed flights. Show the first 5 row.

```
flights %>% arrange(desc(dep_delay)) %>% head(5)
```

```
## # A tibble: 5 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>
                                              900
## 1
      2013
                1
                      9
                              641
                                                        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
                                             1845
                             1139
                                                        1014
                                                                  1457
                                                                                  2210
## 5
      2013
                7
                     22
                              845
                                             1600
                                                        1005
                                                                  1044
                                                                                  1815
## # 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 <dttm>
```

6. Generate a scatter plot with x-axis dist and y-axis delay, where each dot is a unique flights and destination, dist is the average distance of each destination dest, and delay is the average delay time arr\_delay, with the size of dot equals to the count of delay records.

```
flights %>%
group_by(flight, dest) %>%
summarise(delay = mean(arr_delay), dist = mean(distance), n = n()) %>%
ggplot() +
geom_point(aes(x = dist, y = delay, size = n))
## `summarise()` has grouped output by 'flight'. You can override using the
```

```
## `.groups` argument.
```



## Warning: Removed 2824 rows containing missing values or values outside the scale range
## (`geom\_point()`).