# Homework #2: Transforming Data

## Aditi Madhok

## 2/3/21

Let's add our libraries first:

```
library(tidyverse)
library(nycflights13)
```

1. Consider the `flights` variable from the `nycflights13` package. Use the `select` command to create tibbles with the variables described below:

a. Only the carrier and tail number.

```
planes<- select(flights, carrier | tailnum)
planes
```

```
## # A tibble: 336,776 x 2
##    carrier tailnum
##    <chr>   <chr>
##  1 UA      N14228
##  2 UA      N24211
##  3 AA      N619AA
##  4 B6      N804JB
##  5 DL      N668DN
##  6 UA      N39463
##  7 B6      N516JB
##  8 EV      N829AS
##  9 B6      N593JB
## 10 AA      N3ALAA
## # ... with 336,766 more rows
```

b. All variables except the year.

```
not_year<- select (flights, !(year))
not_year
```

```
## # A tibble: 336,776 x 18
##    month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int>    <int>          <int>     <dbl>    <int>          <int>
##  1     1     1      517            515         2      830            819
##  2     1     1      533            529         4      850            830
##  3     1     1      542            540         2      923            850
```

```
## 4      1     1        544              545         -1     1004             1022
## 5      1     1        554              600         -6      812              837
## 6      1     1        554              558         -4      740              728
## 7      1     1        555              600         -5      913              854
## 8      1     1        557              600         -3      709              723
## 9      1     1        557              600         -3      838              846
## 10     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 <dttm>
```

c. Any variable ending with the word 'time'.

```
time_variables<- select (flights, ends_with("time"))
time_variables
```

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

d. The first 9 variables.

```
first_nine<- select(flights, year : arr_delay)
first_nine
```

```
## # A tibble: 336,776 x 9
##      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 1 more variable: arr_delay <dbl>
```

2. Use the filter function to find all the flights that satisfy the following conditions.

2

a. Had an arrival delay of two or more hours.

```
arr_delay_ge2 <- filter(flights, arr_delay>=120)
arr_delay_ge2
```

```
## # 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
## # ... with 10,190 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 <dttm>
```

b. Flew to Houston (airport codes 'IAH' or 'HOU').

```
to_houston <- filter(flights,dest =="IAH" | dest == "HOU")
to_houston
```

```
## # 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 <dttm>
```

c. Departed from JFK in July.

```
JFK_July<-filter(flights,origin =="JFK", month == "7")
JFK_July
```

```
## # A tibble: 10,023 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       44           2150        174      300           100
## 5   2013     7     1       46           2051        235      304          2358
## 6   2013     7     1       48           2001        287      308          2305
## 7   2013     7     1       58           2155        183      335            43
## 8   2013     7     1      100           2146        194      327            30
## 9   2013     7     1      100           2245        135      337           135
## 10  2013     7     1      107           2245        142      158          2359
## # ... with 10,013 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 <dttm>
```

d. Another useful dplyr filtering helper is `between`. Look up what it does and how to use it. Then, use it to find flights that left between 0 and 60 minutes late.

```
late <- filter(flights, (between(dep_delay,0,60)))
late
```

```
## # A tibble: 118,365 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      559            559         0      702            706
## 5   2013     1     1      600            600         0      851            858
## 6   2013     1     1      600            600         0      837            825
## 7   2013     1     1      601            600         1      844            850
## 8   2013     1     1      607            607         0      858            915
## 9   2013     1     1      608            600         8      807            735
## 10  2013     1     1      611            600        11      945            931
## # ... with 118,355 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 <dttm>
```

e. Filter the 'flights' dataset to *remove* all flights with missing departure times.

```
flights_with_dep_time<-filter(flights,!is.na(dep_time))
flights_with_dep_time
```

```
## # A tibble: 328,521 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 328,511 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 <dttm>
```

---

3. Practice with mutate.

   a. Consider the `distance` variable in the `flights` dataset. Currently this is measured in miles. Convert
      this to feet with the `mutate` command (the convereted variable should still be called 'distance').

```
flights_with_feet <- mutate(flights, distance=distance*5280)
flights_with_feet
```

```
## # 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 <dttm>
```

   b. Add a variable `speed` to the `flights` table that gives the average flight speed, **in miles per hour**.

```
flights_with_speed<- mutate(flights,speed=distance/(air_time/60))
flights_with_speed
```

```
## # A tibble: 336,776 x 20
##      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 12 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>,
## #   speed <dbl>
```

    c. Add a variable to `flights` called `early` which is TRUE if the flight arrival early and FALSE if it arrived on time or late.

```
early_flights<- mutate(flights,early=(arr_delay<0))
early_flights
```

```
## # A tibble: 336,776 x 20
##     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 12 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>,
## #   early <lgl>
```

---

4. The `arrange` function sorts a variable from low to high.

    a. Sort 'flights' so that the flights that departed closest to their scheduled departure time are first.

```
most_on_time<-arrange(flights, abs(dep_delay))
most_on_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      559            559         0      702            706
## 2   2013     1     1      600            600         0      851            858
## 3   2013     1     1      600            600         0      837            825
## 4   2013     1     1      607            607         0      858            915
## 5   2013     1     1      615            615         0     1039           1100
## 6   2013     1     1      615            615         0      833            842
## 7   2013     1     1      635            635         0     1028            940
## 8   2013     1     1      655            655         0     1021           1030
```

```
## 9   2013     1     1      739             739        0      1104             1038
## 10  2013     1     1      745             745        0      1135             1125
## # ... 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 <dttm>
```

b. Sort `flights` according to their arrival delay, **from high to low**.

```
latest_arrival<- arrange(flights,desc(arr_delay))
latest_arrival
```

```
## # 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 <dttm>
```

c. Use `arrange` to sort `early_flights` (from Problem 3c) on the variable `early`. (Is `TRUE` or `FALSE` the lower value?) #false is lower value so

```
sorted_by_boolean<-arrange(early_flights)
sorted_by_boolean
```

```
## # A tibble: 336,776 x 20
##       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 12 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>,
## #   early <lgl>
```