# DATA SCIENCE AND ANALYTICS
## Week 2

Ayşe Ceren Çiçek

## Data Transformation

Most of the code examples written down below, taken from this website https://r4ds.had.co.nz/transform.html

## TABLE OF CONTENTS

## Import libraries

```r
#install.packages("nycflights13")
library(nycflights13)
library(tidyverse)
```

## Dataset

The dataset contains data for all flights that departed NYC (i.e. JFK, LGA or EWR) in 2013. Column informations are listed below:

- year, month, day: Date of departure.

- dep_time, arr_time: Actual departure and arrival times (format HHMM or HMM), local tz.

- sched_dep_time, sched_arr_time: Scheduled departure and arrival times (format HHMM or HMM), local tz.

- dep_delay, arr_delay: Departure and arrival delays, in minutes. Negative times represent early departures/arrivals.

- carrier: Two letter carrier abbreviation. See airlines to get name.

- flight: Flight number.

- tailnum: Plane tail number. See planes for additional metadata.

- origin, dest: Origin and destination. See airports for additional metadata.

- air_time: Amount of time spent in the air, in minutes.

- distance: Distance between airports, in miles.

- hour, minute: Time of scheduled departure broken into hour and minutes.

- time_hour: Scheduled date and hour of the flight as a POSIXct date. Along with origin, can be used to join flights data to weather data.

```
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 <dttm>
```

To see the whole dataset, you can run View(flights).

```
summary(flights)
```

```
##       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
## 1st Qu.: 8.00   1st Qu.:2013-04-04 13:00:00
## Median :29.00   Median :2013-07-03 10:00:00
## Mean   :26.23   Mean   :2013-07-03 05:22:54
## 3rd Qu.:44.00   3rd Qu.:2013-10-01 07:00:00
## Max.   :59.00   Max.   :2013-12-31 23:00:00
##
```

## dplyr Basics

- Pick observations by their values (filter()).
- Reorder the rows (arrange()).
- Pick variables by their names (select()).
- Create new variables with functions of existing variables (mutate()).
- Collapse many values down to a single summary (summa rize()).

### Filter

filter() allows you to subset observations based on their values.

We can select all flights on May 1st with.

```
filter(flights, month == 5, day == 1)
```

```
## # A tibble: 964 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     5     1       9          1655        434     308          2020
## 2  2013     5     1     451           500         -9     641           640
## 3  2013     5     1     537           540         -3     836           840
## 4  2013     5     1     544           545         -1     818           827
## 5  2013     5     1     548           600        -12     831           854
## 6  2013     5     1     549           600        -11     804           810
## 7  2013     5     1     553           600         -7     700           712
```

```
## 8  2013     5     1      553          600       -7      655          701
## 9  2013     5     1      554          600       -6      731          756
## 10 2013     5     1      554          600       -6      707          725
## # ... with 954 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>
```

```
(june25 <- filter(flights, month == 6, day == 25))
```

```
## # A tibble: 993 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     6    25        1           2130       151      249             14
## 2  2013     6    25        7           2130       157      237           2359
## 3  2013     6    25       11           2245        86      137              3
## 4  2013     6    25       12           2250        82      143             14
## 5  2013     6    25       27           2146       161      307             30
## 6  2013     6    25       27           2359        28      411            350
## 7  2013     6    25       32           2231       121      409            226
## 8  2013     6    25      103           2359        64      431            344
## 9  2013     6    25      104           1900       364      319           2147
## 10 2013     6    25      118           2300       138      207              8
## # ... with 983 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>
```

```
filter(flights, month == 6 | month == 12)
```

**Logical Operators**

```
## # A tibble: 56,378 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     1       13           2359        14      446            445
## 2  2013    12     1       17           2359        18      443            437
## 3  2013    12     1      453            500        -7      636            651
## 4  2013    12     1      520            515         5      749            808
## 5  2013    12     1      536            540        -4      845            850
## 6  2013    12     1      540            550       -10     1005           1027
## 7  2013    12     1      541            545        -4      734            755
## 8  2013    12     1      546            545         1      826            835
## 9  2013    12     1      549            600       -11      648            659
## 10 2013    12     1      550            600       -10      825            854
## # ... with 56,368 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>
```

x %in% y -> this will select every row where x is one of the values in y

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

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

```
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 <dttm>
```

```
x <- NA
is.na(x)
```

**Missing Values**

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

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

```
filter(df, is.na(x) | x > 1)
```

```
## # A tibble: 2 x 1
##        x
##    <dbl>
## 1     NA
## 2      3
```

**Arrange**

arrange() takes a data frame and a set of column names to order by.

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

```
## # 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>
```

desc() to re-order by a column in descending order.

```
arrange(flights, desc(dep_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     6    27      959           1900       899     1236           2226
## 9   2013     7    22     2257            759       898      121           1026
## 10  2013    12     5      756           1700       896     1058           2020
## # ... 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>
```

**Select**

select() allows you to rapidly zoom in on a useful subset using operations based on the names of the variables.

```
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 <dttm>
```

```r
#everything() moves some variables to the start of the data frame
select(flights, time_hour, air_time, day, everything())
```

```
## # A tibble: 336,776 x 19
##    time_hour           air_time   day  year month dep_time sched_dep_time
##    <dttm>                 <dbl> <int> <int> <int>    <int>          <int>
## 1 2013-01-01 05:00:00      227     1  2013     1      517            515
## 2 2013-01-01 05:00:00      227     1  2013     1      533            529
## 3 2013-01-01 05:00:00      160     1  2013     1      542            540
## 4 2013-01-01 05:00:00      183     1  2013     1      544            545
## 5 2013-01-01 06:00:00      116     1  2013     1      554            600
## 6 2013-01-01 05:00:00      150     1  2013     1      554            558
## 7 2013-01-01 06:00:00      158     1  2013     1      555            600
## 8 2013-01-01 06:00:00       53     1  2013     1      557            600
## 9 2013-01-01 06:00:00      140     1  2013     1      557            600
## 10 2013-01-01 06:00:00     138     1  2013     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>
```

**Mutate**

mutate() always adds new columns at the end of the dataset.

```r
flights_sml <- select(flights,
  year:day,
  ends_with("delay"),
  distance,
  air_time
)

mutate(flights_sml,
  gain = dep_delay - arr_delay,
  speed = distance / air_time * 60
)
```

```
## # A tibble: 336,776 x 9
##    year month   day dep_delay arr_delay distance air_time  gain speed
##   <int> <int> <int>     <dbl>     <dbl>    <dbl>    <dbl> <dbl> <dbl>
## 1  2013     1     1         2        11     1400      227    -9  370.
## 2  2013     1     1         4        20     1416      227   -16  374.
## 3  2013     1     1         2        33     1089      160   -31  408.
## 4  2013     1     1        -1       -18     1576      183    17  517.
## 5  2013     1     1        -6       -25      762      116    19  394.
```

```
##  6  2013     1     1        -4        12       719       150   -16  288.
##  7  2013     1     1        -5        19      1065       158   -24  404.
##  8  2013     1     1        -3       -14       229        53    11  259.
##  9  2013     1     1        -3        -8       944       140     5  405.
## 10  2013     1     1        -2         8       733       138   -10  319.
## # ... with 336,766 more rows
```

If you only want to keep the new variables, use transmute().

```
transmute(flights,
  gain = dep_delay - arr_delay,
  hours = air_time / 60,
  gain_per_hour = gain / hours
)
```

```
## # A tibble: 336,776 x 3
##      gain hours gain_per_hour
##     <dbl> <dbl>         <dbl>
##  1    -9 3.78          -2.38
##  2   -16 3.78          -4.23
##  3   -31 2.67         -11.6
##  4    17 3.05           5.57
##  5    19 1.93           9.83
##  6   -16 2.5           -6.4
##  7   -24 2.63          -9.11
##  8    11 0.883         12.5
##  9     5 2.33           2.14
## 10   -10 2.3           -4.35
## # ... with 336,766 more rows
```

```
transmute(flights,
  dep_time,
  hour = dep_time %/% 100,
  minute = dep_time %% 100
)
```

```
## # A tibble: 336,776 x 3
##    dep_time  hour minute
##       <int> <dbl>  <dbl>
##  1      517     5     17
##  2      533     5     33
##  3      542     5     42
##  4      544     5     44
##  5      554     5     54
##  6      554     5     54
##  7      555     5     55
##  8      557     5     57
##  9      557     5     57
## 10      558     5     58
## # ... with 336,766 more rows
```

```r
(x <- 1:10)
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10
```

```r
lag(x)
```

```
##  [1] NA  1  2  3  4  5  6  7  8  9
```

```r
lead(x)
```

```
##  [1]  2  3  4  5  6  7  8  9 10 NA
```

```r
cumsum(x)
```

```
##  [1]  1  3  6 10 15 21 28 36 45 55
```

```r
#cumulative means
cummean(x)
```

```
##  [1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0 5.5
```

**Summarise**

summarise() collapses a data frame to a single row.

```r
summarise(flights, delay = mean(dep_delay, na.rm = TRUE))
```

```
## # A tibble: 1 x 1
##   delay
##   <dbl>
## 1  12.6
```

```r
by_day <- group_by(flights, year, month, day)
summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##     year month   day delay
##    <int> <int> <int> <dbl>
##  1  2013     1     1 11.5
##  2  2013     1     2 13.9
##  3  2013     1     3 11.0
##  4  2013     1     4  8.95
##  5  2013     1     5  5.73
##  6  2013     1     6  7.15
##  7  2013     1     7  5.42
##  8  2013     1     8  2.55
##  9  2013     1     9  2.28
## 10  2013     1    10  2.84
## # ... with 355 more rows
```

```
flights %>%
  group_by(year, month, day) %>%
  summarise(mean = mean(dep_delay, na.rm = TRUE))
```

**Combining multiple operations with the pipe(%>%)**

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##     year month   day  mean
##    <int> <int> <int> <dbl>
##  1  2013     1     1 11.5
##  2  2013     1     2 13.9
##  3  2013     1     3 11.0
##  4  2013     1     4  8.95
##  5  2013     1     5  5.73
##  6  2013     1     6  7.15
##  7  2013     1     7  5.42
##  8  2013     1     8  2.55
##  9  2013     1     9  2.28
## 10  2013     1    10  2.84
## # ... with 355 more rows
```

- Counts
- Useful Summary Functions
- Grouping by multiple variables
- Ungrouping