# Transforming Data in R

*Pramod Duvvuri*

*3/29/2019*

## Data Manipulation in R

```
library(nycflights13) # Data to be used
library(tidyverse)

## -- Attaching packages ----------------------------------------------------------- tidyverse 1.2.1 --

## v ggplot2 3.1.0       v purrr   0.3.2
## v tibble  2.1.1       v dplyr   0.8.0.1
## v tidyr   0.8.3       v stringr 1.4.0
## v readr   1.3.1       v forcats 0.4.0

## -- Conflicts ------------------------------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

nycflights13::flights

## # A tibble: 336,776 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time
##    <int> <int> <int>   <int>          <int>     <dbl>   <int>
## 1   2013     1     1     517            515         2     830
## 2   2013     1     1     533            529         4     850
## 3   2013     1     1     542            540         2     923
## 4   2013     1     1     544            545        -1    1004
## 5   2013     1     1     554            600        -6     812
## 6   2013     1     1     554            558        -4     740
## 7   2013     1     1     555            600        -5     913
## 8   2013     1     1     557            600        -3     709
## 9   2013     1     1     557            600        -3     838
## 10  2013     1     1     558            600        -2     753
## # ... with 336,766 more rows, and 12 more variables: sched_arr_time <int>,
## #   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>

attach(flights)
?flights
```

### Basics of dplyr

The below functions will be covered in this file.

```
1. filter()
2. arrange()
3. select()
4. mutate()
```

```
5. summarise()
6. group_by()
```

**filter()**

```r
# Select all flights on January 1st
filter(flights, month == 1, day == 1)
```

```
## # A tibble: 842 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>
## 1   2013     1     1      517            515         2      830
## 2   2013     1     1      533            529         4      850
## 3   2013     1     1      542            540         2      923
## 4   2013     1     1      544            545        -1     1004
## 5   2013     1     1      554            600        -6      812
## 6   2013     1     1      554            558        -4      740
## 7   2013     1     1      555            600        -5      913
## 8   2013     1     1      557            600        -3      709
## 9   2013     1     1      557            600        -3      838
## 10  2013     1     1      558            600        -2      753
## # ... with 832 more rows, and 12 more variables: sched_arr_time <int>,
## #   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>
```

```r
jan_1_data <- filter(flights, month == 1, day == 1)
# To assing the data to another variable and print to console
(dec25_data <- filter(flights, month == 12, day == 25))
```

```
## # A tibble: 719 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>
## 1   2013    12    25      456            500        -4      649
## 2   2013    12    25      524            515         9      805
## 3   2013    12    25      542            540         2      832
## 4   2013    12    25      546            550        -4     1022
## 5   2013    12    25      556            600        -4      730
## 6   2013    12    25      557            600        -3      743
## 7   2013    12    25      557            600        -3      818
## 8   2013    12    25      559            600        -1      855
## 9   2013    12    25      559            600        -1      849
## 10  2013    12    25      600            600         0      850
## # ... with 709 more rows, and 12 more variables: sched_arr_time <int>,
## #   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>
```