# Data transformation

## Yifan Jin

## 09/02/2020

## Introduction

Visualisation is an important tool for insight generation, but it is rare that you get the data in exactly the right form you need. We may need to create, summarises, rename, reorder vairables.

### Prerequisites

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

```
## -- Attaching packages -------------- tidyverse 1.3.0 --
```

```
## v ggplot2 3.2.1     v purrr   0.3.3
## v tibble  2.1.3     v dplyr   0.8.3
## v tidyr   1.0.0     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 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>
```

`int`: intergers

`dbl`: doubles, or real numbers

`chr`: character vectors, or strings

`dttm`: date-times (a date+ a time)

`lgl`: logical, vectors that contain only `TRUE` or `FALSE`

`fctr`: factors, which R uses to represent categorical variables with fixed possible values

`date`: dates

**dplyr basics**

`filter()`: pick observations by their values

`arrange()`: reorder the rows

`select()`: pick variables by their names

`mutate()`: create new variables with functions of existing variable

`summarise()`: collapse many values down to single summary

All verbs work similarly:

1. The first argument is a data frame.

2. The subsequent arguments describe what to do with the data frame, using the variable names (without quotes)

3. The result is a new data frame

## Filter rows with filter()

`filter()` allows you to subset observations based on their values. The first argument is the name of the data frame. The second and subsequent arguments are the expressions that filter the data frame.

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

```
## # A tibble: 842 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 832 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>
```

```r
jan1<-filter(flights,month==1,day==1)

(dec25<-filter(flights,month==12,day==25))
```

```
## # A tibble: 719 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    25      456            500        -4      649            651
## 2   2013    12    25      524            515         9      805            814
## 3   2013    12    25      542            540         2      832            850
## 4   2013    12    25      546            550        -4     1022           1027
## 5   2013    12    25      556            600        -4      730            745
## 6   2013    12    25      557            600        -3      743            752
## 7   2013    12    25      557            600        -3      818            831
## 8   2013    12    25      559            600        -1      855            856
## 9   2013    12    25      559            600        -1      849            855
## 10  2013    12    25      600            600         0      850            846
## # ... with 709 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>
```
```r
# brackets that we can print the result
```

### Comparisons

Comparison operator: >, >=, <,<=, != and ==

```r
# filter(flights,month=1) This is false since we need to use ==
sqrt(2)^2==2
```

```
## [1] FALSE
```
```r
1/49*49==1
```

```
## [1] FALSE
```
```r
# This result is surprising, when we compare two numbers, one is floating, one is integer

#Instead of relying on `==`, we use near()

near(sqrt(2)^2,2)
```

```
## [1] TRUE
```
```r
near(1/49*49,1)
```

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

### Logical operators

&: and

|: or

!: not

```r
filter(flights,month==11|month==12)
```

```
## # 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>
```

```r
# or

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

**Missing values**

```r
NA>5
```

```
## [1] NA
```

```r
10==NA
```

```
## [1] NA
```

```r
NA+10
```

```
## [1] NA
```

```r
NA/2
```

```
## [1] NA
```

```r
NA==NA
```

```
## [1] NA
```

```r
X<-NA
Y<-NA
X==Y
```

```
## [1] NA
```

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

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

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

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

**Exercises**

**1.Find all flights that had an arrival delay of two or more hours**

**Answer**

```
filter(flights,arr_delay>120)
```

```
## # A tibble: 10,034 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 10,024 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>
```

**2. Find all flights that Flew to Houston (IAH or HOU)**

**Answer**

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

```
## # 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>
```

**3. Find all flights that were operated by United, American, or Delta**

**Answer**

```
filter(flights,carrier=="UA"|carrier=="AA"|carrier=="DL")
```

```
## # A tibble: 139,504 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      554            600        -6      812            837
## 5   2013     1     1      554            558        -4      740            728
## 6   2013     1     1      558            600        -2      753            745
## 7   2013     1     1      558            600        -2      924            917
## 8   2013     1     1      558            600        -2      923            937
## 9   2013     1     1      559            600        -1      941            910
## 10  2013     1     1      559            600        -1      854            902
## # ... with 139,494 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>
```

**4. Departed in summer (July, August, and September)**

**Answer**

```
filter(flights,month==7|month==8|month==9)
```

```
## # A tibble: 86,326 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       43           2130       193      322             14
## 5   2013     7     1       44           2150       174      300            100
## 6   2013     7     1       46           2051       235      304           2358
## 7   2013     7     1       48           2001       287      308           2305
## 8   2013     7     1       58           2155       183      335             43
## 9   2013     7     1      100           2146       194      327             30
## 10  2013     7     1      100           2245       135      337            135
## # ... with 86,316 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>
```

**5. Arrived more than two hours late, but didn't leave late**

```
filter(flights,arr_delay>120|dep_delay<=0)
```

```
## # A tibble: 210,094 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      544            545        -1     1004           1022
## 2   2013     1     1      554            600        -6      812            837
## 3   2013     1     1      554            558        -4      740            728
## 4   2013     1     1      555            600        -5      913            854
## 5   2013     1     1      557            600        -3      709            723
```

```
## 6   2013     1     1      557          600         -3      838             846
## 7   2013     1     1      558          600         -2      753             745
## 8   2013     1     1      558          600         -2      849             851
## 9   2013     1     1      558          600         -2      853             856
## 10  2013     1     1      558          600         -2      924             917
## # ... with 210,084 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>
```

**6.Were delayed by at least an hours late, but made up over 30 minutes in flight**

**Answer**

```r
filter(flights,arr_delay<dep_delay-30|dep_delay>=60)
```

```
## # A tibble: 43,165 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      701          700          1     1123           1154
## 2   2013     1     1      811          630        101     1047            830
## 3   2013     1     1      820          820          0     1249           1329
## 4   2013     1     1      826          715         71     1136           1045
## 5   2013     1     1      840          845         -5     1311           1350
## 6   2013     1     1      848         1835        853     1001           1950
## 7   2013     1     1      857          851          6     1157           1222
## 8   2013     1     1      909          810         59     1331           1315
## 9   2013     1     1      957          733        144     1056            853
## 10  2013     1     1     1025          951         34     1258           1302
## # ... with 43,155 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>
```

**7. Departed between midnight and 6am (inclusive)**

**Answer**

```r
filter(flights,dep_time>=000&dep_time<=600)
```

```
## # A tibble: 9,344 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 9,334 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>
```

**8. Another useful dplyr filtering helper is `between()`. What does it do?**

**Answer**

```
filter(flights,between(dep_time,000,600))
```

```
## # A tibble: 9,344 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 9,334 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,is.na(dep_time))
```

```
## # A tibble: 8,255 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       NA           1630        NA       NA           1815
## 2   2013     1     1       NA           1935        NA       NA           2240
## 3   2013     1     1       NA           1500        NA       NA           1825
## 4   2013     1     1       NA            600        NA       NA            901
## 5   2013     1     2       NA           1540        NA       NA           1747
## 6   2013     1     2       NA           1620        NA       NA           1746
## 7   2013     1     2       NA           1355        NA       NA           1459
## 8   2013     1     2       NA           1420        NA       NA           1644
## 9   2013     1     2       NA           1321        NA       NA           1536
## 10  2013     1     2       NA           1545        NA       NA           1910
## # ... with 8,245 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>
```

## Arrange rows with arrange()