

NYC Flights 2013

John Cruz

2023-02-11

Required Libraries

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

Using dplyr

Row Operations

filter()

```
# Flights that departed on January 1
flights |>
  filter(month == 1 & day == 1)
```

```
## # A tibble: 842 x 19
##   year month   day dep_time sched_de~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier
##   <int> <int> <int>   <int>      <int>   <dbl>   <int>   <int>   <dbl> <chr>
## 1  2013     1     1     517         515     2     830     819     11 UA
## 2  2013     1     1     533         529     4     850     830     20 UA
## 3  2013     1     1     542         540     2     923     850     33 AA
## 4  2013     1     1     544         545    -1    1004    1022    -18 B6
## 5  2013     1     1     554         600    -6     812     837    -25 DL
## 6  2013     1     1     554         558    -4     740     728     12 UA
## 7  2013     1     1     555         600    -5     913     854     19 B6
## 8  2013     1     1     557         600    -3     709     723    -14 EV
## 9  2013     1     1     557         600    -3     838     846     -8 B6
## 10 2013     1     1     558         600    -2     753     745      8 AA
## # ... with 832 more rows, 9 more variables: flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>, and abbreviated variable names
## #   1: sched_dep_time, 2: dep_delay, 3: arr_time, 4: sched_arr_time,
## #   5: arr_delay
```

```
# Flights that departed in January or February
flights |>
  filter(month %in% c(1, 2))
```

```
## # A tibble: 51,955 x 19
##   year month   day dep_time sched_de~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier
##   <int> <int> <int>   <int>      <int>    <dbl>   <int>    <int>    <dbl> <chr>
## 1  2013     1     1     517        515      2     830     819      11 UA
## 2  2013     1     1     533        529      4     850     830     20 UA
## 3  2013     1     1     542        540      2     923     850     33 AA
## 4  2013     1     1     544        545     -1    1004    1022    -18 B6
## 5  2013     1     1     554        600     -6     812     837    -25 DL
## 6  2013     1     1     554        558     -4     740     728     12 UA
## 7  2013     1     1     555        600     -5     913     854     19 B6
## 8  2013     1     1     557        600     -3     709     723    -14 EV
## 9  2013     1     1     557        600     -3     838     846     -8 B6
##10  2013     1     1     558        600     -2     753     745      8 AA
## # ... with 51,945 more rows, 9 more variables: flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>, and abbreviated variable names
## #   1: sched_dep_time, 2: dep_delay, 3: arr_time, 4: sched_arr_time,
## #   5: arr_delay
```

arrange()

```
flights |>
  arrange(desc(dep_delay))
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_de~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier
##   <int> <int> <int>   <int>      <int>    <dbl>   <int>    <int>    <dbl> <chr>
## 1  2013     1     9     641        900    1301    1242    1530    1272 HA
## 2  2013     6    15    1432       1935    1137    1607    2120    1127 MQ
## 3  2013     1    10    1121       1635    1126    1239    1810    1109 MQ
## 4  2013     9    20    1139       1845    1014    1457    2210    1007 AA
## 5  2013     7    22     845       1600    1005    1044    1815     989 MQ
## 6  2013     4    10    1100       1900     960    1342    2211     931 DL
## 7  2013     3    17    2321        810     911     135    1020     915 DL
## 8  2013     6    27     959       1900     899    1236    2226     850 DL
## 9  2013     7    22    2257        759     898     121    1026     895 DL
##10  2013    12     5     756       1700     896    1058    2020     878 AA
## # ... with 336,766 more rows, 9 more variables: flight <int>, tailnum <chr>,
## #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
## #   minute <dbl>, time_hour <dtm>, and abbreviated variable names
## #   1: sched_dep_time, 2: dep_delay, 3: arr_time, 4: sched_arr_time,
## #   5: arr_delay
```

distinct()

```
flights |>
  distinct(origin, dest)
```

```
## # A tibble: 224 x 2
##   origin dest
##   <chr> <chr>
## 1 EWR    IAH
## 2 LGA    IAH
## 3 JFK    MIA
## 4 JFK    BQN
## 5 LGA    ATL
## 6 EWR    ORD
## 7 EWR    FLL
## 8 LGA    IAD
## 9 JFK    MCO
## 10 LGA    ORD
## # ... with 214 more rows
```

Column Operations

mutate()

- *.before* or *.after* “Determine new columns placement in data frame.”
- *.keep* “Control which variables are kept. (‘used’ argument keeps the inputs from your calculations)”

```
flights |>
  mutate(
    gain = dep_delay - arr_delay,
    hours = air_time / 60,
    gain_per_hour = gain / hours,
    .keep = "used"
  )
```

```
## # A tibble: 336,776 x 6
##   dep_delay arr_delay air_time gain hours gain_per_hour
##   <dbl>    <dbl>    <dbl> <dbl> <dbl>    <dbl>
## 1         2        11      227    -9  3.78     -2.38
## 2         4        20      227   -16  3.78     -4.23
## 3         2        33      160   -31  2.67    -11.6
## 4        -1       -18      183    17  3.05      5.57
## 5        -6       -25      116    19  1.93      9.83
## 6        -4        12      150   -16  2.5     -6.4
## 7        -5        19      158   -24  2.63    -9.11
## 8        -3       -14       53    11  0.883    12.5
## 9        -3        -8      140     5  2.33      2.14
## 10       -2         8      138   -10  2.3     -4.35
## # ... with 336,766 more rows
```

select()

```
flights |>
  mutate(
    gain = dep_delay - arr_delay,
    hours = air_time / 60,
    gain_per_hour = gain / hours,
    .keep = "used"
  )
```

```
## # A tibble: 336,776 x 6
##   dep_delay arr_delay air_time  gain hours gain_per_hour
##   <dbl>     <dbl>    <dbl> <dbl> <dbl>         <dbl>
## 1         2         11     227    -9  3.78         -2.38
## 2         4         20     227   -16  3.78         -4.23
## 3         2         33     160   -31  2.67        -11.6
## 4        -1        -18     183    17  3.05          5.57
## 5        -6        -25     116    19  1.93          9.83
## 6        -4         12     150   -16  2.5          -6.4
## 7        -5         19     158   -24  2.63         -9.11
## 8        -3        -14      53    11  0.883         12.5
## 9        -3         -8     140     5  2.33          2.14
## 10       -2          8     138   -10  2.3         -4.35
## # ... with 336,766 more rows
```