# Data Manipulation Verbs

## Ibrahim

## 15/03/2022

Here we apply the data manipulation verbs on `nycflight13` dataset to handle this data and extract pieces of information.

```r
# call the required libraries
library(tidyverse)
```

**Filter**

```
## Warning: package 'tidyverse' was built under R version 4.0.5
```

```
## -- Attaching packages --------------------------------------- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5     v purrr   0.3.4
## v tibble  3.1.1     v dplyr   1.0.6
## v tidyr   1.1.3     v stringr 1.4.0
## v readr   1.4.0     v forcats 0.5.1
```

```
## Warning: package 'ggplot2' was built under R version 4.0.5
```

```
## Warning: package 'tibble' was built under R version 4.0.5
```

```
## Warning: package 'tidyr' was built under R version 4.0.5
```

```
## Warning: package 'readr' was built under R version 4.0.5
```

```
## Warning: package 'dplyr' was built under R version 4.0.5
```

```
## Warning: package 'forcats' was built under R version 4.0.5
```

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

```r
library(nycflights13)
```

```
## Warning: package 'nycflights13' was built under R version 4.0.5
```

```r
df <- flights
# filter based on conditions
filter(df, 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
filter(df, 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
# you may combine conditions using logical operators
filter(df, month == 1 | month == 12)
```

```
## # A tibble: 55,139 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 55,129 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 by combining variables in vector
filter(df, month %in% c(11, 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
# comma means and (&)
# for example, extract all records for flights, that were not delayed (arr and dep) more than 2 hrs
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>
```

```r
# call the required libraries
library(tidyverse)
```

```
library(nycflights13)
df <- flights
# Select columns by name
select(flights, year, month, day)
```

**Select**

```
## # 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 (inclusive)
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 (inclusive)
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
# rename() is a variant of select() that keeps all the variables that aren't explicitly mentioned:
rename(flights, tail_num = tailnum)
```

```
## # 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>, tail_num <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

```r
# Move a variable to the start of the data frame.
select(flights, time_hour, air_time, everything())
```

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

```r
# call the required libraries
library(tidyverse)
```

```
library(nycflights13)
df <- flights
# Create a new dataset
flights_sml <- select(df, year:day, ends_with("delay"), distance, air_time )
# add new columns to the data frame
mutate(flights_sml, gain = dep_delay - arr_delay, speed = distance / air_time * 60)
```

**Mutate**

```
## # 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
```

```
# Note that you can refer to columns that you've just created:
mutate(flights_sml, gain = dep_delay - arr_delay, hours = air_time / 60, gain_per_hour = gain / hours )
```

```
## # A tibble: 336,776 x 10
##     year month   day dep_delay arr_delay distance air_time  gain hours
##    <int> <int> <int>    <dbl>     <dbl>    <dbl>    <dbl> <dbl> <dbl>
##  1  2013     1     1        2        11     1400      227    -9 3.78
##  2  2013     1     1        4        20     1416      227   -16 3.78
##  3  2013     1     1        2        33     1089      160   -31 2.67
##  4  2013     1     1       -1       -18     1576      183    17 3.05
##  5  2013     1     1       -6       -25      762      116    19 1.93
##  6  2013     1     1       -4        12      719      150   -16 2.5
##  7  2013     1     1       -5        19     1065      158   -24 2.63
##  8  2013     1     1       -3       -14      229       53    11 0.883
##  9  2013     1     1       -3        -8      944      140     5 2.33
## 10  2013     1     1       -2         8      733      138   -10 2.3
## # ... with 336,766 more rows, and 1 more variable: gain_per_hour <dbl>
```

```
# 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
```

```
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# sort data by distance
arrange(df, distance)
```

**Arrange**

```
## # 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     7    27       NA            106        NA       NA            245
## 2    2013     1     3     2127           2129        -2     2222           2224
## 3    2013     1     4     1240           1200        40     1333           1306
## 4    2013     1     4     1829           1615       134     1937           1721
## 5    2013     1     4     2128           2129        -1     2218           2224
## 6    2013     1     5     1155           1200        -5     1241           1306
## 7    2013     1     6     2125           2129        -4     2224           2224
## 8    2013     1     7     2124           2129        -5     2212           2224
## 9    2013     1     8     2127           2130        -3     2304           2225
## 10   2013     1     9     2126           2129        -3     2217           2224
## # ... 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>
```

```
# sort data by distance descendingly
arrange(df, desc(distance))
```

```
## # 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      857            900        -3     1516           1530
## 2    2013     1     2      909            900         9     1525           1530
## 3    2013     1     3      914            900        14     1504           1530
## 4    2013     1     4      900            900         0     1516           1530
## 5    2013     1     5      858            900        -2     1519           1530
## 6    2013     1     6     1019            900        79     1558           1530
## 7    2013     1     7     1042            900       102     1620           1530
## 8    2013     1     8      901            900         1     1504           1530
## 9    2013     1     9      641            900      1301     1242           1530
## 10   2013     1    10      859            900        -1     1449           1530
```

```
## # ... 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>
```

```r
# Sort Data by Multiple Variables
arrange(df, dep_time, arr_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     6    24        1           1950       251      105           2130
## 2   2013     4    10        1           1930       271      106           2101
## 3   2013     1    13        1           2249        72      108           2357
## 4   2013     2    11        1           2100       181      111           2225
## 5   2013     3    19        1           2250        71      120              5
## 6   2013     2    24        1           2245        76      121           2354
## 7   2013     1    31        1           2100       181      124           2225
## 8   2013     7    22        1           2305        56      135             13
## 9   2013     5    22        1           1935       266      154           2140
## 10  2013     7     1        1           2029       212      236           2359
## # ... 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>
```

```r
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# extract a statistical metric from variable / variables of the data
summarise(df, delay = mean(dep_delay, na.rm = TRUE))
```

**Summarise**

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

```r
# group the data of the flights by the date
by_day <- group_by(flights, year, month, day)
# get the average delay per date/day
summarise(by_day, delay = mean(dep_delay, na.rm = TRUE)) # Imagine that we want to explore the relation
```

**Grouping**

```
## `summarise()` has grouped output by 'year', 'month'. You can override using the '.groups' argument.
```

```
## # 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
```

```r
by_dest <- group_by(flights, dest)

# extract the number of flights, average distance and average delay for each destination
delay <- summarise(by_dest, count= n(), dist= mean(distance, na.rm = TRUE), delay= mean(arr_delay, na.r

# visualise to understand the relationship
ggplot(data= delay, mapping=aes(x= dist, y= delay)) + geom_point(aes(size= count), alpha= 1/ 3) + geom_
```
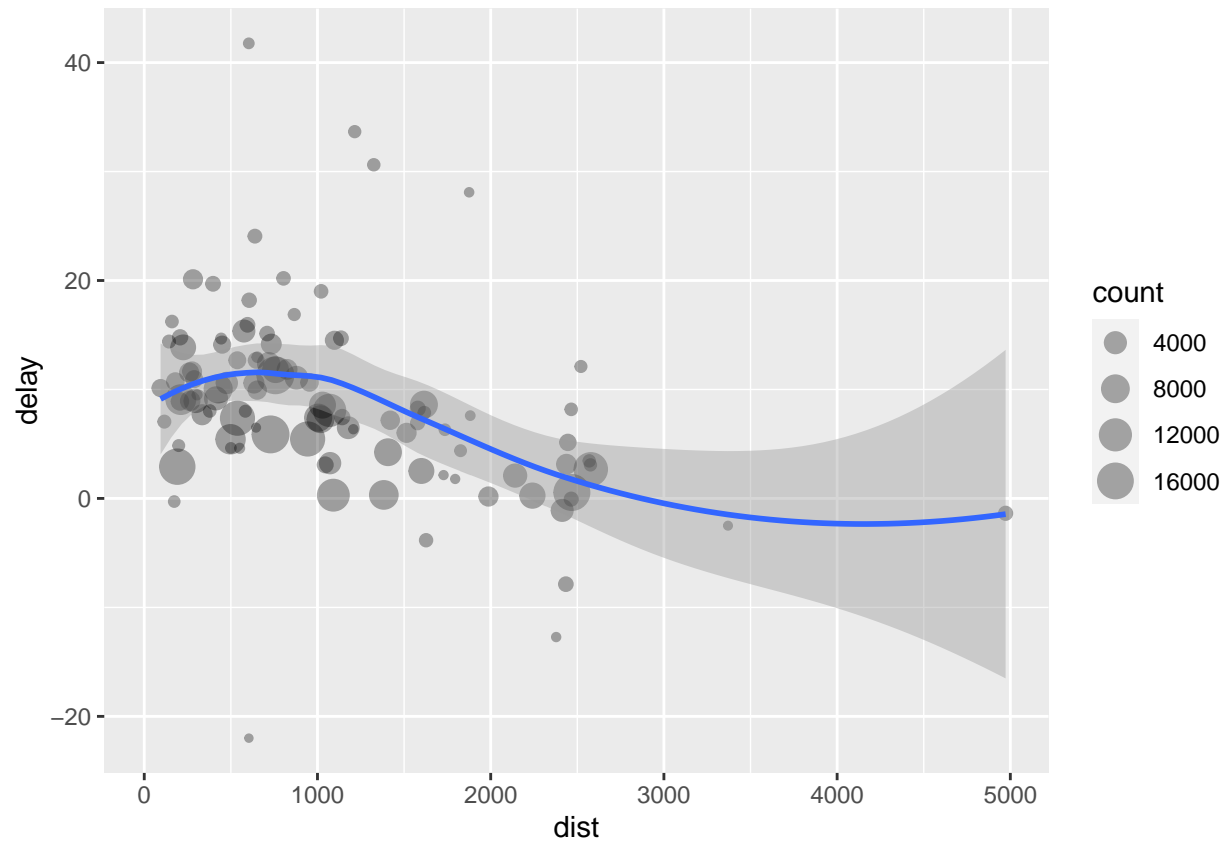
```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

```
## Warning: Removed 1 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
df %>%
group_by(dest) %>%
summarise(count= n(), dist= mean(distance, na.rm = TRUE), delay= mean(arr_delay, na.rm = TRUE)) %>%
filter(count > 20, dest != 'HNL') %>%
ggplot(mapping=aes(x= dist, y= delay)) +
  geom_point(aes(size= count), alpha= 1/ 3) +
  geom_smooth()
```

**Pipe Operator**

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```