# dplyr Basics

Sadaf Zuhra

2024-09-07

## dplyr Basics

The five key dplyr functions that allow you to solve the vast majority of your data-manipulation challenges:

- 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 (summarize()).

These can all be used in conjunction with group_by(), which changes the scope of each function from operating on the entire dataset to operating on it group-by-group. ## Filter Rows with filter()

```
library(tidyverse)
library(nycflights13)
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
## # i 832 more rows
## # i 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>
```

```
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
## # i 709 more rows
## # i 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>
```

```
library(tidyverse)
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
## # i 55,393 more rows
## # i 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>
```

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

```
library(tidyverse)
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
## # i 316,040 more rows
## # i 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(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
## # i 316,040 more rows
## # i 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()

```r
library(tidyverse)
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
## # i 336,766 more rows
## # i 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(flights, desc(arr_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     7    22     2257            759       898      121           1026
## 9   2013    12     5      756           1700       896     1058           2020
## 10  2013     5     3     1133           2055       878     1250           2215
## # i 336,766 more rows
## # i 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 Columns with select()

```
library(tidyverse)
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
## # i 336,766 more rows
```

```
library(tidyverse)
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
## # i 336,766 more rows
```

```
library(tidyverse)
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
## # i 336,766 more rows
## # i 9 more variables: flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

```
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
## # i 336,766 more rows
## # i 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>
```

```
library(tidyverse)
select(flights, contains("TIME"))
```

```
## # A tibble: 336,776 x 6
```

```
##    dep_time sched_dep_time arr_time sched_arr_time air_time time_hour
##       <int>          <int>    <int>          <int>    <dbl> <dttm>
## 1      517            515      830            819      227 2013-01-01 05:00:00
## 2      533            529      850            830      227 2013-01-01 05:00:00
## 3      542            540      923            850      160 2013-01-01 05:00:00
## 4      544            545     1004           1022      183 2013-01-01 05:00:00
## 5      554            600      812            837      116 2013-01-01 06:00:00
## 6      554            558      740            728      150 2013-01-01 05:00:00
## 7      555            600      913            854      158 2013-01-01 06:00:00
## 8      557            600      709            723       53 2013-01-01 06:00:00
## 9      557            600      838            846      140 2013-01-01 06:00:00
## 10     558            600      753            745      138 2013-01-01 06:00:00
## # i 336,766 more rows
```

**change the name of column with rename()**

```
library(tidyverse)
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
## # i 336,766 more rows
## # i 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>
```

## Add New Variables with mutate()

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

```
library(tidyverse)
mutate(flights_sml, gain = arr_delay - dep_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.
## # i 336,766 more rows
```

```
library(tidyverse)
mutate(flights_sml, gain = arr_delay - dep_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
## # i 336,766 more rows
## # i 1 more variable: gain_per_hour <dbl>
```

```
library(tidyverse)
transmute(flights, gain = arr_delay - dep_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
## # i 336,766 more rows
```

**Grouped Summaries with summarize()**

```
library(tidyverse)
summarize(flights, delay = mean(dep_delay, na.rm = TRUE))
```

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

summarize() is not terribly useful unless we pair it with group_by().

```
library(tidyverse)
by_day<- group_by(flights, year, month, day)
summarize(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
## # i 355 more rows
```