# Chapter 3: Data Transformation in R

Francesco Calzona

## SECTION 3.1

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

```
-- Attaching core tidyverse packages ---------------------- tidyverse 2.0.0 --
v dplyr     1.1.4     v readr     2.1.5
v forcats   1.0.0     v stringr   1.5.1
v ggplot2   3.4.4     v tibble    3.2.1
v lubridate 1.9.3     v tidyr     1.3.1
v purrr     1.0.2
-- Conflicts --------------------------------------- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to becon
```

```r
View(flights) #View is useful to scroll datasets; otherwise print/glimpse
glimpse(flights)
```

```
Rows: 336,776
Columns: 19
$ year          <int> 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2~
$ month         <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ day           <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ dep_time      <int> 517, 533, 542, 544, 554, 554, 555, 557, 557, 558, 558, ~
$ sched_dep_time <int> 515, 529, 540, 545, 600, 558, 600, 600, 600, 600, 600, ~
$ dep_delay     <dbl> 2, 4, 2, -1, -6, -4, -5, -3, -3, -2, -2, -2, -2, -2, -1~
$ arr_time      <int> 830, 850, 923, 1004, 812, 740, 913, 709, 838, 753, 849,~
$ sched_arr_time <int> 819, 830, 850, 1022, 837, 728, 854, 723, 846, 745, 851,~
```

```
$ arr_delay    <dbl> 11, 20, 33, -18, -25, 12, 19, -14, -8, 8, -2, -3, 7, -1~
$ carrier      <chr> "UA", "UA", "AA", "B6", "DL", "UA", "B6", "EV", "B6", "~
$ flight       <int> 1545, 1714, 1141, 725, 461, 1696, 507, 5708, 79, 301, 4~
$ tailnum      <chr> "N14228", "N24211", "N619AA", "N804JB", "N668DN", "N394~
$ origin       <chr> "EWR", "LGA", "JFK", "JFK", "LGA", "EWR", "EWR", "LGA",~
$ dest         <chr> "IAH", "IAH", "MIA", "BQN", "ATL", "ORD", "FLL", "IAD",~
$ air_time     <dbl> 227, 227, 160, 183, 116, 150, 158, 53, 140, 138, 149, 1~
$ distance     <dbl> 1400, 1416, 1089, 1576, 762, 719, 1065, 229, 944, 733, ~
$ hour         <dbl> 5, 5, 5, 5, 6, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 5, 6, 6, 6~
$ minute       <dbl> 15, 29, 40, 45, 0, 58, 0, 0, 0, 0, 0, 0, 0, 0, 0, 59, 0~
$ time_hour    <dttm> 2013-01-01 05:00:00, 2013-01-01 05:00:00, 2013-01-01 0~
```

A glimpse into pipes:

```
flights |>
  filter(dest == "IAH") |>
  group_by(year, month, day) |>
  summarize(arr_delay = mean(arr_delay, na.rm = TRUE))
```

`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 arr_delay
   <int> <int> <int>     <dbl>
 1  2013     1     1     17.8
 2  2013     1     2      7
 3  2013     1     3     18.3
 4  2013     1     4     -3.2
 5  2013     1     5     20.2
 6  2013     1     6      9.28
 7  2013     1     7     -7.74
 8  2013     1     8      7.79
 9  2013     1     9     18.1
10  2013     1    10      6.68
# i 355 more rows
```

## Row operations

```
flights |> filter(dep_delay > 120)
```

```
# A tibble: 9,723 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     848           1835       853     1001           1950
 2  2013     1     1     957            733       144     1056            853
 3  2013     1     1    1114            900       134     1447           1222
 4  2013     1     1    1540           1338       122     2020           1825
 5  2013     1     1    1815           1325       290     2120           1542
 6  2013     1     1    1842           1422       260     1958           1535
 7  2013     1     1    1856           1645       131     2212           2005
 8  2013     1     1    1934           1725       129     2126           1855
 9  2013     1     1    1938           1703       155     2109           1823
10  2013     1     1    1942           1705       157     2124           1830
# i 9,713 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>
```

The %in% operator combines == and |:

```
flights |> filter(day %in% c(1,3,5))
```

```
# A tibble: 33,105 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 33,095 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>
```

dplyr never modifies the input: if we want to save the result, it must be assigned to a new variable. Arrange orders stuff in an increasing order.

```
flights |> arrange(year, month, day, dep_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     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>
```

We can also do the opposite:

```
flights |> arrange(desc(year), desc(month), desc(day), desc(dep_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    12    31     2356           2359        -3      436            445
 2  2013    12    31     2355           2359        -4      430            440
 3  2013    12    31     2332           2245        47       58              3
 4  2013    12    31     2328           2330        -2      412            409
 5  2013    12    31     2321           2250        31       46              8
 6  2013    12    31     2310           2255        15        7           2356
```

```
 7  2013    12     31     2245            2250            -5      2359            2356
 8  2013    12     31     2235            2245           -10      2351            2355
 9  2013    12     31     2218            2219            -1       315             304
10  2013    12     31     2211            2159            12       100              45
# 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>
```

distinct() finds all unique rows, on all or some columns:

```
flights|> distinct()
```

```
# 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>
```

```
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
# i 214 more rows
```

```
  flights |> distinct(origin, dest, .keep_all = TRUE)
```

```
# A tibble: 224 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 214 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>
```

count also gives you the number of occurrances.

```
  flights |> count(origin, dest, sort = TRUE)
```

```
# A tibble: 224 x 3
   origin dest      n
   <chr>  <chr> <int>
 1 JFK    LAX   11262
 2 LGA    ATL   10263
 3 LGA    ORD    8857
 4 JFK    SFO    8204
 5 LGA    CLT    6168
```

```
 6 EWR    ORD    6100
 7 JFK    BOS    5898
 8 LGA    MIA    5781
 9 JFK    MCO    5464
10 EWR    BOS    5327
# i 214 more rows
```

**Exercises 3.2.5**

```
glimpse(flights)
```

```
Rows: 336,776
Columns: 19
$ year          <int> 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2~
$ month         <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ day           <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ dep_time      <int> 517, 533, 542, 544, 554, 554, 555, 557, 557, 558, 558, ~
$ sched_dep_time <int> 515, 529, 540, 545, 600, 558, 600, 600, 600, 600, 600, ~
$ dep_delay     <dbl> 2, 4, 2, -1, -6, -4, -5, -3, -3, -2, -2, -2, -2, -2, -1~
$ arr_time      <int> 830, 850, 923, 1004, 812, 740, 913, 709, 838, 753, 849,~
$ sched_arr_time <int> 819, 830, 850, 1022, 837, 728, 854, 723, 846, 745, 851,~
$ arr_delay     <dbl> 11, 20, 33, -18, -25, 12, 19, -14, -8, 8, -2, -3, 7, -1~
$ carrier       <chr> "UA", "UA", "AA", "B6", "DL", "UA", "B6", "EV", "B6", "~
$ flight        <int> 1545, 1714, 1141, 725, 461, 1696, 507, 5708, 79, 301, 4~
$ tailnum       <chr> "N14228", "N24211", "N619AA", "N804JB", "N668DN", "N394~
$ origin        <chr> "EWR", "LGA", "JFK", "JFK", "LGA", "EWR", "EWR", "LGA",~
$ dest          <chr> "IAH", "IAH", "MIA", "BQN", "ATL", "ORD", "FLL", "IAD",~
$ air_time      <dbl> 227, 227, 160, 183, 116, 150, 158, 53, 140, 138, 149, 1~
$ distance      <dbl> 1400, 1416, 1089, 1576, 762, 719, 1065, 229, 944, 733, ~
$ hour          <dbl> 5, 5, 5, 5, 6, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 5, 6, 6, 6~
$ minute        <dbl> 15, 29, 40, 45, 0, 58, 0, 0, 0, 0, 0, 0, 0, 0, 0, 59, 0~
$ time_hour     <dttm> 2013-01-01 05:00:00, 2013-01-01 05:00:00, 2013-01-01 0~
```

```
flights |> filter(arr_delay >= 120)
```

```
# A tibble: 10,200 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
# i 10,190 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>
```

```
  flights |> filter(dest %in% c("IAH", "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
# i 9,303 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>
```

```
  flights |> filter(carrier %in% c("UA", "AA", "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
# i 139,494 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>
```

```
  flights |> filter(month %in% c(7, 8, 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
# i 86,316 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>
```

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

```
# A tibble: 3 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    10     7     1350           1350         0     1736           1526
```

```
2  2013     5    23    1810          1810           0    2208          2000
3  2013     7     1     905           905           0    1443          1223
# 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
flights |> filter(dep_delay >= 60 & arr_delay <= dep_delay - 30)
```

```
# A tibble: 2,074 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     1716           1545        91     2140           2039
 2  2013     1     1     2205           1720       285       46           2040
 3  2013     1     1     2326           2130       116      131             18
 4  2013     1     3     1503           1221       162     1803           1555
 5  2013     1     3     1821           1530       171     2131           1910
 6  2013     1     3     1839           1700        99     2056           1950
 7  2013     1     3     1850           1745        65     2148           2120
 8  2013     1     3     1923           1815        68     2036           1958
 9  2013     1     3     1941           1759       102     2246           2139
10  2013     1     3     1950           1845        65     2228           2227
# i 2,064 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
flights |> arrange(desc(dep_delay), dep_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     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     6    27      959           1900       899     1236           2226
 9  2013     7    22     2257            759       898      121           1026
```

```
10  2013    12     5       756           1700          896     1058            2020
```
```
# 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>
```

```r
flights |> arrange(desc(distance/air_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     5    25     1709           1700         9     1923           1937
 2  2013     7     2     1558           1513        45     1745           1719
 3  2013     5    13     2040           2025        15     2225           2226
 4  2013     3    23     1914           1910         4     2045           2043
 5  2013     1    12     1559           1600        -1     1849           1917
 6  2013    11    17      650            655        -5     1059           1150
 7  2013     2    21     2355           2358        -3      412            438
 8  2013    11    17      759            800        -1     1212           1255
 9  2013    11    16     2003           1925        38       17             36
10  2013    11    16     2349           2359       -10      402            440
# 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>
```

```r
flights |> filter(year == 2013) |> count(month, day)
```

```
# A tibble: 365 x 3
   month   day     n
   <int> <int> <int>
 1     1     1   842
 2     1     2   943
 3     1     3   914
 4     1     4   915
 5     1     5   720
 6     1     6   832
 7     1     7   933
 8     1     8   899
 9     1     9   902
```

```
10      1    10    932
# i 355 more rows
```

```
  flights |> filter(distance == max(distance))
```

```
# A tibble: 342 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
# i 332 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>
```

```
  flights |> filter(distance == min(distance))
```

```
# A tibble: 1 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
# 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>
```

## Column operations

Mutate is used to create new columns from existing ones. It is possible to add them .before.

```
  flights |>
    mutate(
```

```
    gain = dep_delay - arr_delay,
    speed = distance / air_time * 60,
    .before = 3
  )
```

```
# A tibble: 336,776 x 21
    year month  gain speed   day dep_time sched_dep_time dep_delay arr_time
   <int> <int> <dbl> <dbl> <int>    <int>          <int>     <dbl>    <int>
 1  2013     1    -9  370.     1      517            515         2      830
 2  2013     1   -16  374.     1      533            529         4      850
 3  2013     1   -31  408.     1      542            540         2      923
 4  2013     1    17  517.     1      544            545        -1     1004
 5  2013     1    19  394.     1      554            600        -6      812
 6  2013     1   -16  288.     1      554            558        -4      740
 7  2013     1   -24  404.     1      555            600        -5      913
 8  2013     1    11  259.     1      557            600        -3      709
 9  2013     1     5  405.     1      557            600        -3      838
10  2013     1   -10  319.     1      558            600        -2      753
# i 336,766 more rows
# i 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>
```

```
  flights |>
    mutate(
      gain = dep_delay - arr_delay,
      speed = distance / air_time * 60,
      .after = dep_time
    )
```

```
# A tibble: 336,776 x 21
    year month   day dep_time  gain speed sched_dep_time dep_delay arr_time
   <int> <int> <int>    <int> <dbl> <dbl>          <int>     <dbl>    <int>
 1  2013     1     1      517    -9  370.            515         2      830
 2  2013     1     1      533   -16  374.            529         4      850
 3  2013     1     1      542   -31  408.            540         2      923
 4  2013     1     1      544    17  517.            545        -1     1004
 5  2013     1     1      554    19  394.            600        -6      812
 6  2013     1     1      554   -16  288.            558        -4      740
 7  2013     1     1      555   -24  404.            600        -5      913
```

```
 8  2013      1      1       557    11  259.               600        -3       709
 9  2013      1      1       557     5  405.               600        -3       838
10  2013      1      1       558   -10  319.               600        -2       753
# i 336,766 more rows
# i 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>
```

It is possible to keep only the columns used during the calculations.

```
flights |>
  mutate(
    gain = dep_delay - arr_delay,
    speed = distance / air_time,
    .keep = "used"
  )
```

```
# A tibble: 336,776 x 6
   dep_delay arr_delay air_time distance  gain speed
       <dbl>     <dbl>    <dbl>    <dbl> <dbl> <dbl>
 1         2        11      227     1400    -9  6.17
 2         4        20      227     1416   -16  6.24
 3         2        33      160     1089   -31  6.81
 4        -1       -18      183     1576    17  8.61
 5        -6       -25      116      762    19  6.57
 6        -4        12      150      719   -16  4.79
 7        -5        19      158     1065   -24  6.74
 8        -3       -14       53      229    11  4.32
 9        -3        -8      140      944     5  6.74
10        -2         8      138      733   -10  5.31
# i 336,766 more rows
```

select() is for selecting columns:

```
flights |> select(air_time, distance)
```

```
# A tibble: 336,776 x 2
  air_time distance
     <dbl>    <dbl>
1      227     1400
2      227     1416
```

```
 3       160    1089
 4       183    1576
 5       116     762
 6       150     719
 7       158    1065
 8        53     229
 9       140     944
10       138     733
# i 336,766 more rows
```

```
flights |> select(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
```

```
flights |> select(!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>
```

```
  flights |> select(where(is.character))
```

```
# A tibble: 336,776 x 4
   carrier tailnum origin dest
   <chr>   <chr>   <chr>  <chr>
 1 UA      N14228  EWR    IAH
 2 UA      N24211  LGA    IAH
 3 AA      N619AA  JFK    MIA
 4 B6      N804JB  JFK    BQN
 5 DL      N668DN  LGA    ATL
 6 UA      N39463  EWR    ORD
 7 B6      N516JB  EWR    FLL
 8 EV      N829AS  LGA    IAD
 9 B6      N593JB  JFK    MCO
10 AA      N3ALAA  LGA    ORD
# i 336,766 more rows
```

```
  flights |> select(starts_with("dep"))
```

```
# A tibble: 336,776 x 2
   dep_time dep_delay
      <int>     <dbl>
 1      517         2
 2      533         4
 3      542         2
 4      544        -1
 5      554        -6
 6      554        -4
 7      555        -5
 8      557        -3
 9      557        -3
10      558        -2
# i 336,766 more rows
```

```
flights |> select(ends_with("a"))
```

```
# A tibble: 336,776 x 0
```

```
flights |> select(contains("a"))
```

```
# A tibble: 336,776 x 10
     year   day dep_delay arr_time sched_arr_time arr_delay carrier tailnum
    <int> <int>     <dbl>    <int>          <int>     <dbl> <chr>   <chr>
 1   2013     1         2      830            819        11 UA      N14228
 2   2013     1         4      850            830        20 UA      N24211
 3   2013     1         2      923            850        33 AA      N619AA
 4   2013     1        -1     1004           1022       -18 B6      N804JB
 5   2013     1        -6      812            837       -25 DL      N668DN
 6   2013     1        -4      740            728        12 UA      N39463
 7   2013     1        -5      913            854        19 B6      N516JB
 8   2013     1        -3      709            723       -14 EV      N829AS
 9   2013     1        -3      838            846        -8 B6      N593JB
10   2013     1        -2      753            745         8 AA      N3ALAA
# i 336,766 more rows
# i 2 more variables: air_time <dbl>, distance <dbl>
```

```
flights |> select(num_range("x", 1:3))
```

```
# A tibble: 336,776 x 0
```

There are even more in ?select.

```
flights |> select(tail_num = tailnum)
```

```
# A tibble: 336,776 x 1
   tail_num
   <chr>
 1 N14228
 2 N24211
 3 N619AA
 4 N804JB
```

```
 5 N668DN
 6 N39463
 7 N516JB
 8 N829AS
 9 N593JB
10 N3ALAA
# i 336,766 more rows
```

Rename limits ourselves to the last operation, on all columns

```
flights |> rename(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>
```

relocate() changes column positions. By default it takes them to the front, otherwise we
can use .before and .after.

```
flights |> relocate(time_hour, air_time)
```

```
# 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>
```

**Exercises 3.3.5**

```
flights |> select(dep_time, sched_dep_time, dep_delay)
```

```
# A tibble: 336,776 x 3
   dep_time sched_dep_time dep_delay
      <int>          <int>     <dbl>
 1      517            515         2
 2      533            529         4
 3      542            540         2
 4      544            545        -1
 5      554            600        -6
 6      554            558        -4
 7      555            600        -5
 8      557            600        -3
 9      557            600        -3
10      558            600        -2
# i 336,766 more rows
```

```
flights |> select(dep_time, dep_time)
```

```
# A tibble: 336,776 x 1
   dep_time
      <int>
 1      517
 2      533
```

```
 3         542
 4         544
 5         554
 6         554
 7         555
 8         557
 9         557
10         558
# i 336,766 more rows
```

```r
flights |> select(any_of(c("year", "month", "day", "dep_delay", "arr_delay")))
```

```
# A tibble: 336,776 x 5
    year month   day dep_delay arr_delay
   <int> <int> <int>     <dbl>     <dbl>
 1  2013     1     1         2        11
 2  2013     1     1         4        20
 3  2013     1     1         2        33
 4  2013     1     1        -1       -18
 5  2013     1     1        -6       -25
 6  2013     1     1        -4        12
 7  2013     1     1        -5        19
 8  2013     1     1        -3       -14
 9  2013     1     1        -3        -8
10  2013     1     1        -2         8
# i 336,766 more rows
```

```r
flights |> select(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
```

```
flights |> rename(air_time_min = air_time) |> relocate(air_time_min)
```

```
# A tibble: 336,776 x 19
   air_time_min  year month   day dep_time sched_dep_time dep_delay arr_time
          <dbl> <int> <int> <int>    <int>          <int>     <dbl>    <int>
 1          227  2013     1     1      517            515         2      830
 2          227  2013     1     1      533            529         4      850
 3          160  2013     1     1      542            540         2      923
 4          183  2013     1     1      544            545        -1     1004
 5          116  2013     1     1      554            600        -6      812
 6          150  2013     1     1      554            558        -4      740
 7          158  2013     1     1      555            600        -5      913
 8           53  2013     1     1      557            600        -3      709
 9          140  2013     1     1      557            600        -3      838
10          138  2013     1     1      558            600        -2      753
# i 336,766 more rows
# i 11 more variables: sched_arr_time <int>, arr_delay <dbl>, carrier <chr>,
#   flight <int>, tailnum <chr>, origin <chr>, dest <chr>, distance <dbl>,
#   hour <dbl>, minute <dbl>, time_hour <dttm>
```

## Groups

group_by() divides the dataset into groups.

```
flights |> group_by(month)
```

```
# A tibble: 336,776 x 19
# Groups:   month [12]
   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>
```

The "grouped" feature is referred as **class**.