

# Week6

Nitin

July 23, 2019

```
#Run required librarys.
```

```
library(nycflights13)
```

```
## Warning: package 'nycflights13' was built under R version 3.5.3
```

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.5.3
```

```
## -- Attaching packages ----- tidyverse 1.2.1 --
```

```
## v ggplot2 3.2.0      v purrr  0.2.5
## v tibble  2.1.3      v dplyr  0.8.0.1
## v tidyr   0.8.1      v stringr 1.3.1
## v readr   1.1.1      v forcats 0.3.0
```

```
## Warning: package 'ggplot2' was built under R version 3.5.3
```

```
## Warning: package 'tibble' was built under R version 3.5.3
```

```
## Warning: package 'dplyr' was built under R version 3.5.3
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
```

```
## x dplyr::lag()     masks stats::lag()
```

```
# See flights details
```

```
flights
```

```
## # A tibble: 336,776 x 19
```

```
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     1     517             515           2     830
## 2  2013     1     1     533             529           4     850
## 3  2013     1     1     542             540           2     923
## 4  2013     1     1     544             545          -1    1004
## 5  2013     1     1     554             600          -6     812
## 6  2013     1     1     554             558          -4     740
## 7  2013     1     1     555             600          -5     913
## 8  2013     1     1     557             600          -3     709
## 9  2013     1     1     557             600          -3     838
## 10 2013     1     1     558             600          -2     753
```

```
## # ... with 336,766 more rows, and 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 <dtm>
```

```
#dplyr basics
```

```
#Filter rows with filter()
```

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

```
## # A tibble: 842 x 19
```

```
##   year month   day dep_time sched_dep_time dep_delay arr_time
```

```
##      <int> <int> <int>      <int>      <int>      <dbl>      <int>
## 1  2013      1      1      517      515          2      830
## 2  2013      1      1      533      529          4      850
## 3  2013      1      1      542      540          2      923
## 4  2013      1      1      544      545         -1     1004
## 5  2013      1      1      554      600         -6      812
## 6  2013      1      1      554      558         -4      740
## 7  2013      1      1      555      600         -5      913
## 8  2013      1      1      557      600         -3      709
## 9  2013      1      1      557      600         -3      838
## 10 2013      1      1      558      600         -2      753
## # ... with 832 more rows, and 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 <dtm>
```

```
# If you want to save the result, you'll need to use the assignment operator, <-
jan1 <- filter(flights, month == 1, day == 1)
```

```
# To print out the results and saves them to a variable you can wrap the assignment in parentheses:
(dec25 <- filter(flights, month == 12, day == 25))
```

```
## # A tibble: 719 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>      <int>      <dbl>      <int>
## 1  2013    12    25     456        500        -4       649
## 2  2013    12    25     524        515         9       805
## 3  2013    12    25     542        540         2       832
## 4  2013    12    25     546        550        -4      1022
## 5  2013    12    25     556        600        -4       730
## 6  2013    12    25     557        600        -3       743
## 7  2013    12    25     557        600        -3       818
## 8  2013    12    25     559        600        -1       855
## 9  2013    12    25     559        600        -1       849
## 10 2013    12    25     600        600         0       850
## # ... with 709 more rows, and 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 <dtm>
```

```
# Comparisons
#filter(flights, month = 1)
```

```
# Using ==: floating point numbers
sqrt(2) ^ 2 == 2
```

```
## [1] FALSE
```

```
#Instead of relying on ==, use near():
near(sqrt(2) ^ 2, 2)
```

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

```
near(1 / 49 * 49, 1)
```

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

```
# Logical operators
```

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

```
## # A tibble: 55,403 x 19
```

```
##   year month   day dep_time sched_dep_time dep_delay arr_time
```

```
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
```

```
## 1  2013    11     1       5           2359           6     352
```

```
## 2  2013    11     1      35           2250          105     123
```

```
## 3  2013    11     1     455           500           -5     641
```

```
## 4  2013    11     1     539           545           -6     856
```

```
## 5  2013    11     1     542           545           -3     831
```

```
## 6  2013    11     1     549           600          -11     912
```

```
## 7  2013    11     1     550           600          -10     705
```

```
## 8  2013    11     1     554           600           -6     659
```

```
## 9  2013    11     1     554           600           -6     826
```

```
## 10 2013    11     1     554           600           -6     749
```

```
## # ... with 55,393 more rows, and 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 <dtm>
```

```
#In this short hand x %in% y will select every row where x is one of the values in y
```

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

```
# To find flights that were not delayed by two hours.
```

```
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
```

```
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
```

```
## 1  2013     1     1     517           515           2     830
```

```
## 2  2013     1     1     533           529           4     850
```

```
## 3  2013     1     1     542           540           2     923
```

```
## 4  2013     1     1     544           545          -1    1004
```

```
## 5  2013     1     1     554           600          -6     812
```

```
## 6  2013     1     1     554           558          -4     740
```

```
## 7  2013     1     1     555           600          -5     913
```

```
## 8  2013     1     1     557           600          -3     709
```

```
## 9  2013     1     1     557           600          -3     838
```

```
## 10 2013     1     1     558           600          -2     753
```

```
## # ... with 316,040 more rows, and 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 <dtm>
```

```
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
```

```
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
```

```
## 1  2013     1     1     517           515           2     830
```

```
## 2  2013     1     1     533           529           4     850
```

```
## 3  2013     1     1     542           540           2     923
```

```
## 4  2013     1     1     544           545          -1    1004
```

```
## 5  2013     1     1     554           600          -6     812
```

```
## 6 2013 1 1 554 558 -4 740
## 7 2013 1 1 555 600 -5 913
## 8 2013 1 1 557 600 -3 709
## 9 2013 1 1 557 600 -3 838
## 10 2013 1 1 558 600 -2 753
## # ... with 316,040 more rows, and 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 <dtm>
```

```
#Missing values
```

```
NA > 5
```

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

```
10 == NA
```

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

```
NA + 10
```

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

```
NA / 2
```

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

```
# This is most confusing result
```

```
NA == NA
```

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

```
# Let x be Mary's age. We don't know how old she is.
```

```
x <- NA
```

```
# Let y be John's age. We don't know how old he is.
```

```
y <- NA
```

```
# Are John and Mary the same age?
```

```
x == y
```

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

```
# determine if a value is missing, use is.na()
```

```
is.na(x)
```

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

```
# If want preserve missing values, ask it explicitly
```

```
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
```

Arrange rows with arrange()

```
#Arrange rows with arrange()
arrange(flights, year, month, day)
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     1     517             515           2     830
## 2  2013     1     1     533             529           4     850
## 3  2013     1     1     542             540           2     923
## 4  2013     1     1     544             545          -1    1004
## 5  2013     1     1     554             600          -6     812
## 6  2013     1     1     554             558          -4     740
## 7  2013     1     1     555             600          -5     913
## 8  2013     1     1     557             600          -3     709
## 9  2013     1     1     557             600          -3     838
##10  2013     1     1     558             600          -2     753
## # ... with 336,766 more rows, and 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 <dtm>
```

```
#Missing values are always sorted at the end
arrange(flights, desc(dep_delay))
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     9     641             900        1301    1242
## 2  2013     6    15    1432            1935        1137    1607
## 3  2013     1    10    1121            1635        1126    1239
## 4  2013     9    20    1139            1845        1014    1457
## 5  2013     7    22     845            1600        1005    1044
## 6  2013     4    10    1100            1900         960    1342
## 7  2013     3    17    2321             810         911     135
## 8  2013     6    27     959            1900         899    1236
## 9  2013     7    22    2257             759         898     121
##10  2013    12     5     756            1700         896    1058
## # ... with 336,766 more rows, and 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 <dtm>
```

```
# Creat Df
df <- tibble(x = c(5, 2, NA))
arrange(df, x)
```

```
## # A tibble: 3 x 1
##       x
##   <dbl>
## 1     2
## 2     5
## 3    NA
```

```
arrange(df, desc(x))
```

```
## # A tibble: 3 x 1
##       x
##   <dbl>
## 1     5
## 2     2
## 3    NA
```

```
#Select columns with select()
```

```
# Select columns by name
```

```
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
## # ... 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
##   <int>         <int>         <dbl>   <int>         <int>         <dbl>
## 1     517           515           2     830           819           11
## 2     533           529           4     850           830           20
## 3     542           540           2     923           850           33
## 4     544           545          -1    1004          1022          -18
```

```
## 5      554      600      -6      812      837      -25
## 6      554      558      -4      740      728      12
## 7      555      600      -5      913      854      19
## 8      557      600      -3      709      723      -14
## 9      557      600      -3      838      846      -8
## 10     558      600      -2      753      745      8
## # ... with 336,766 more rows, and 10 more variables: carrier <chr>,
## #   flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
## #   distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
# Use rename(), which is a variant of select()
rename(flights, tail_num = tailnum)
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>
## 1  2013     1     1     517             515           2     830
## 2  2013     1     1     533             529           4     850
## 3  2013     1     1     542             540           2     923
## 4  2013     1     1     544             545          -1    1004
## 5  2013     1     1     554             600          -6     812
## 6  2013     1     1     554             558          -4     740
## 7  2013     1     1     555             600          -5     913
## 8  2013     1     1     557             600          -3     709
## 9  2013     1     1     557             600          -3     838
## 10 2013     1     1     558             600          -2     753
## # ... with 336,766 more rows, and 12 more variables: sched_arr_time <int>,
## #   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 <dtm>
```

```
#Another option is to use select() in along with the everything() helper.
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
##   <dtm>              <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>
```

```
vars <- c("year", "month", "day", "dep_delay", "arr_delay")
```

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

```
## # A tibble: 336,776 x 6
##   dep_time sched_dep_time arr_time sched_arr_time air_time
##   <int>      <int>      <int>      <int>      <dbl>
## 1      517          515      830          819      227
## 2      533          529      850          830      227
## 3      542          540      923          850      160
## 4      544          545     1004         1022      183
## 5      554          600      812          837      116
## 6      554          558      740          728      150
## 7      555          600      913          854      158
## 8      557          600      709          723       53
## 9      557          600      838          846      140
## 10     558          600      753          745      138
## # ... with 336,766 more rows, and 1 more variable: time_hour <dtm>
```

```
#Add new variables using mutate()
flights_sml <- select(flights,
  year:day,
  ends_with("delay"),
  distance,
  air_time
)
mutate(flights_sml,
  gain = dep_delay - arr_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.
## # ... with 336,766 more rows
```

```
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>
```

*#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
```

*# Useful creation functions*

```
transmute(flights,
  dep_time,
  hour = dep_time %/% 100,
  minute = dep_time %% 100
)
```

```
## # A tibble: 336,776 x 3
##   dep_time hour minute
##   <int> <dbl> <dbl>
## 1    517 5 17
## 2    533 5 33
## 3    542 5 42
## 4    544 5 44
## 5    554 5 54
## 6    554 5 54
## 7    555 5 55
## 8    557 5 57
## 9    557 5 57
## 10    558 5 58
## # ... with 336,766 more rows
```

*# Use lead() and lag() allow you to refer to leading or lagging values*  
(x <- 1:10)

```

## [1] 1 2 3 4 5 6 7 8 9 10
lag(x)

## [1] NA 1 2 3 4 5 6 7 8 9
lead(x)

## [1] 2 3 4 5 6 7 8 9 10 NA
# R provides functions for running sums, products, mins and maxes: cumsum(), cumprod(), cummin(), cummax()
x

## [1] 1 2 3 4 5 6 7 8 9 10
cumsum(x)

## [1] 1 3 6 10 15 21 28 36 45 55
cummean(x)

## [1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0 5.5
# Using ranking functions
y <- c(1, 2, 2, NA, 3, 4)
min_rank(y)

## [1] 1 2 2 NA 4 5
min_rank(desc(y))

## [1] 5 3 3 NA 2 1
#Use min_rank() variants row_number(), dense_rank(), percent_rank(), cume_dist()
row_number(y)

## [1] 1 2 3 NA 4 5
dense_rank(y)

## [1] 1 2 2 NA 3 4
percent_rank(y)

## [1] 0.00 0.25 0.25 NA 0.75 1.00
cume_dist(y)

## [1] 0.2 0.6 0.6 NA 0.8 1.0
#Grouped summaries with summarise()
summarise(flights, delay = mean(dep_delay, na.rm = TRUE))

## # A tibble: 1 x 1
##   delay
##   <dbl>
## 1  12.6
#Just summarise() is not terribly useful unless we pair it with group_by()
by_day <- group_by(flights, year, month, day)
summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))

## # A tibble: 365 x 4
## # Groups:   year, month [12]

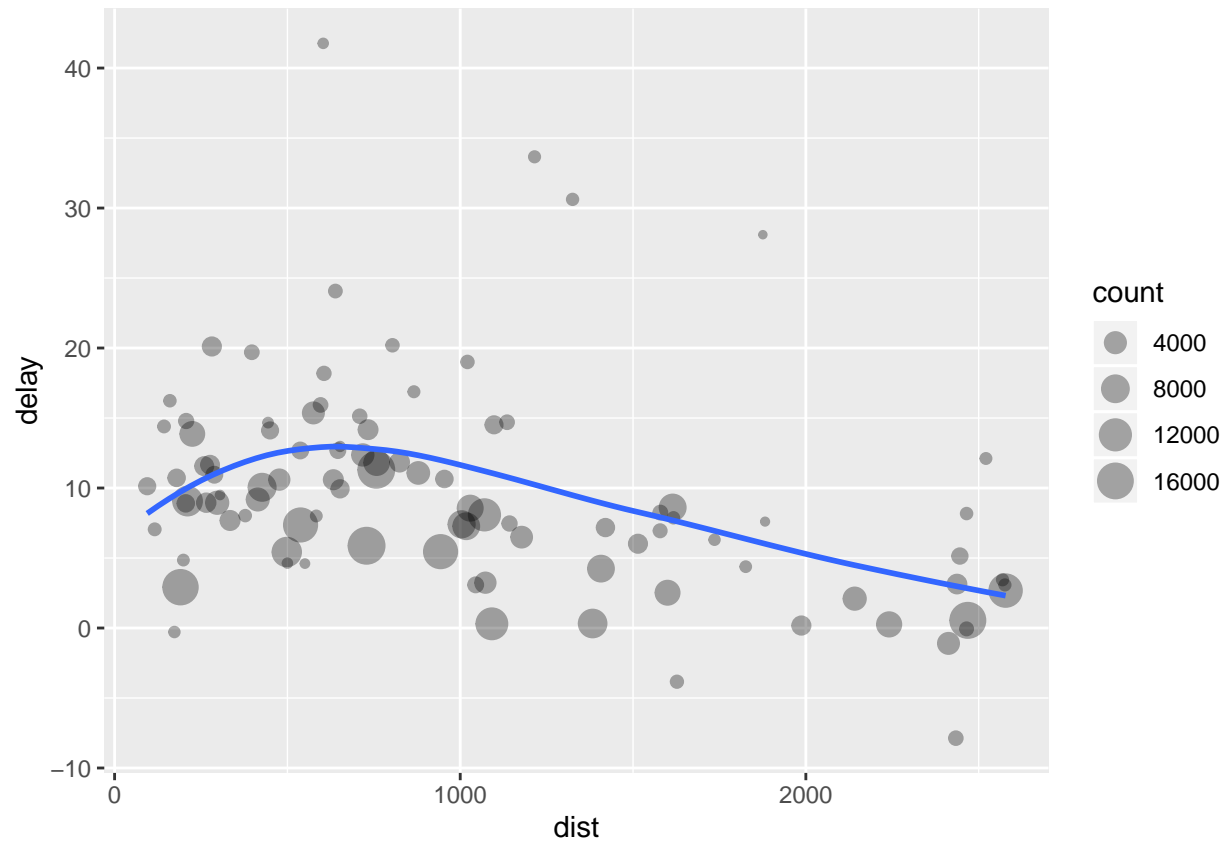
```

```
##      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
## # ... with 355 more rows

#Combining multiple operations with the pipe
by_dest <- group_by(flights, dest)
delay <- summarise(by_dest,
  count = n(),
  dist = mean(distance, na.rm = TRUE),
  delay = mean(arr_delay, na.rm = TRUE)
)
delay <- filter(delay, count > 20, dest != "HNL")

# It looks like delays increase with distance up to ~750 miles
# and then decrease. Maybe as flights get longer there's more
# ability to make up delays in the air?
ggplot(data = delay, mapping = aes(x = dist, y = delay)) +
  geom_point(aes(size = count), alpha = 1/3) +
  geom_smooth(se = FALSE)

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



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

```
#Another way to tackle the problem is to use pipe, %>%
delays <- flights %>%
  group_by(dest) %>%
  summarise(
    count = n(),
    dist = mean(distance, na.rm = TRUE),
    delay = mean(arr_delay, na.rm = TRUE)
  ) %>%
  filter(count > 20, dest != "HNL")
```

Missing values

```
#Missing values
flights %>%
  group_by(year, month, day) %>%
  summarise(mean = mean(dep_delay))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day mean
##   <int> <int> <int> <dbl>
## 1  2013     1     1  NA
## 2  2013     1     2  NA
## 3  2013     1     3  NA
## 4  2013     1     4  NA
```

```
## 5 2013 1 5 NA
## 6 2013 1 6 NA
## 7 2013 1 7 NA
## 8 2013 1 8 NA
## 9 2013 1 9 NA
## 10 2013 1 10 NA
## # ... with 355 more rows
```

*#All aggregation functions have an na.rm argument which removes the missing values before to computation*

```
flights %>%
  group_by(year, month, day) %>%
  summarise(mean = mean(dep_delay, na.rm = TRUE))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day mean
##   <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
```

*#Save this dataset so for reuse*

```
not_cancelled <- flights %>%
  filter(!is.na(dep_delay), !is.na(arr_delay))
```

```
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(mean = mean(dep_delay))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day mean
##   <int> <int> <int> <dbl>
## 1 2013     1     1 11.4
## 2 2013     1     2 13.7
## 3 2013     1     3 10.9
## 4 2013     1     4  8.97
## 5 2013     1     5  5.73
## 6 2013     1     6  7.15
## 7 2013     1     7  5.42
## 8 2013     1     8  2.56
## 9 2013     1     9  2.30
## 10 2013     1    10  2.84
## # ... with 355 more rows
```

*#Counts*

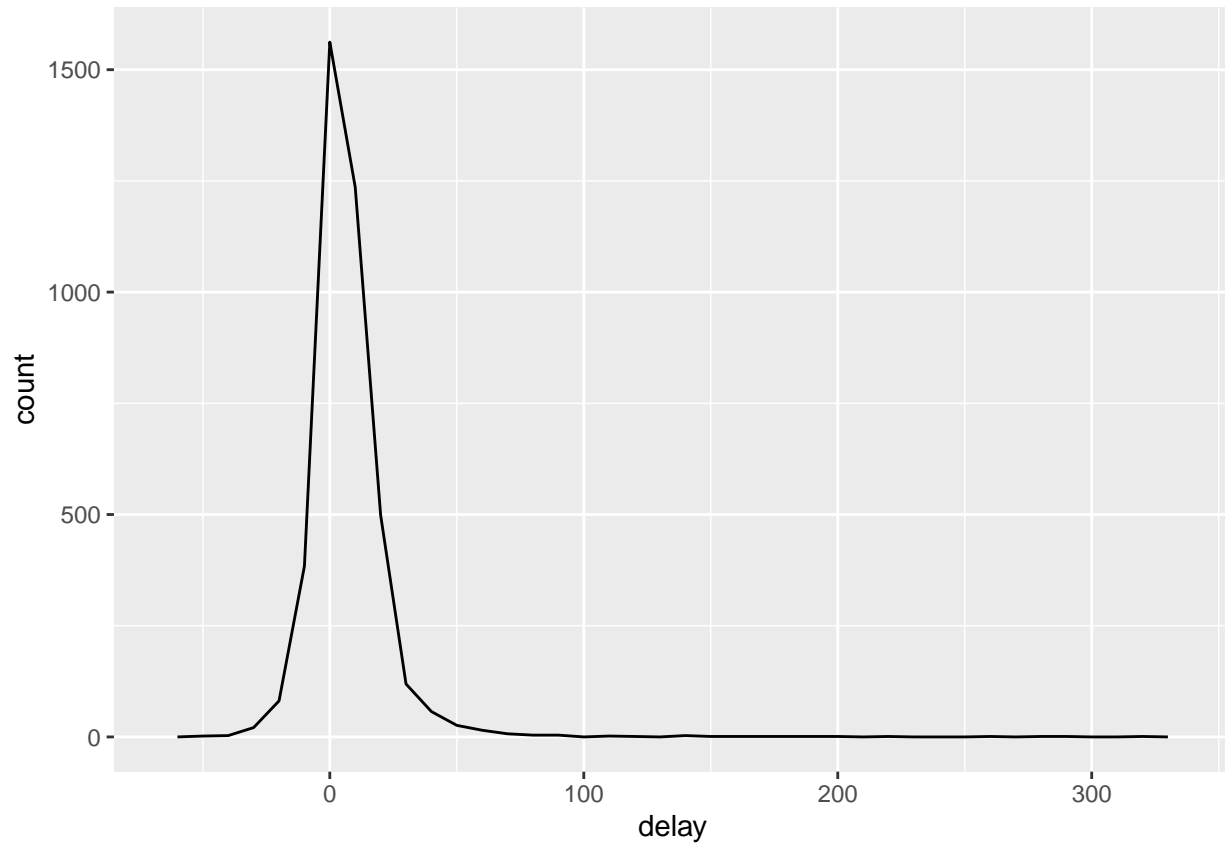
```
delays <- not_cancelled %>%
  group_by(tailnum) %>%
```

```

summarise(
  delay = mean(arr_delay)
)

ggplot(data = delays, mapping = aes(x = delay)) +
  geom_freqpoly(binwidth = 10)

```

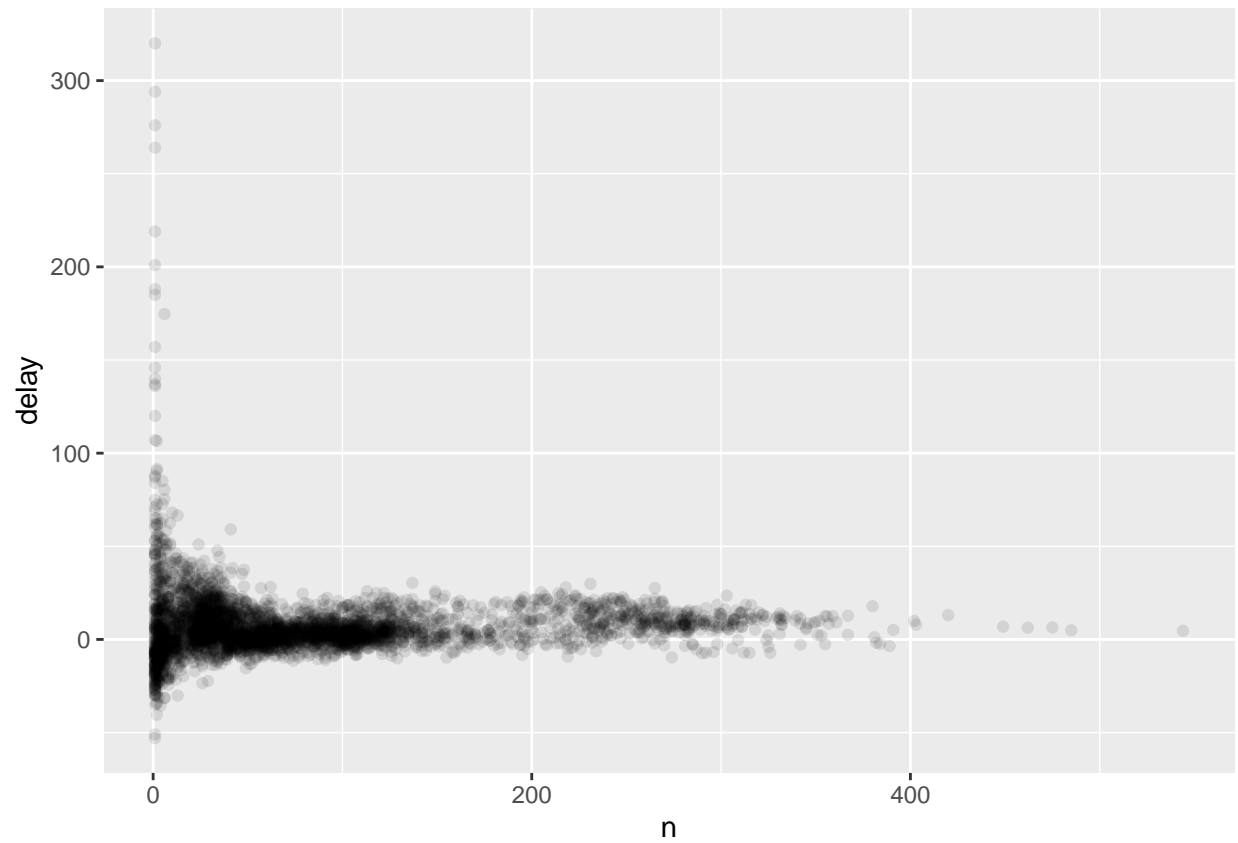


```

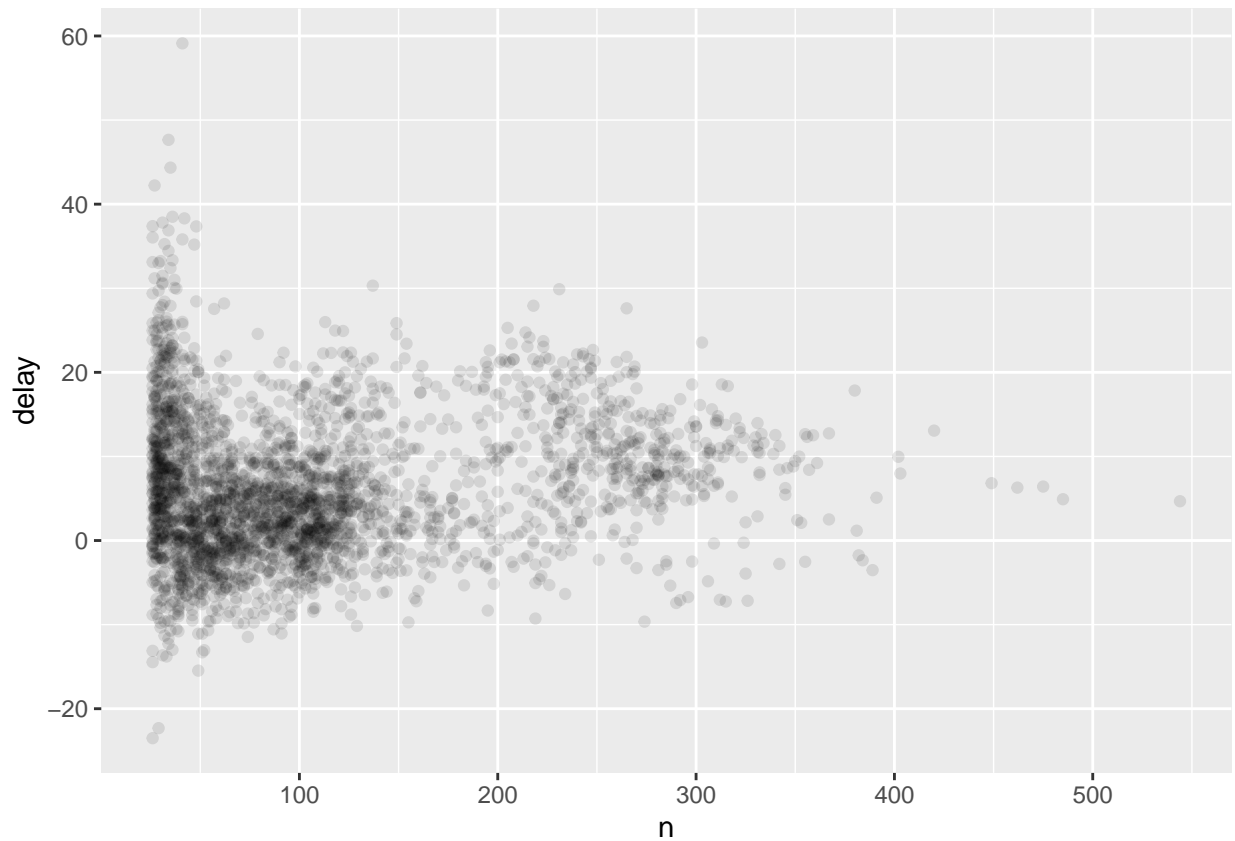
#Get more insight by drawing a scatterplot of number of flights vs. average delay:
delays <- not_cancelled %>%
  group_by(tailnum) %>%
  summarise(
    delay = mean(arr_delay, na.rm = TRUE),
    n = n()
  )

ggplot(data = delays, mapping = aes(x = n, y = delay)) +
  geom_point(alpha = 1/10)

```



```
#It is useful to filter out the groups with the smallest numbers of observations.  
delays %>%  
  filter(n > 25) %>%  
  ggplot(mapping = aes(x = n, y = delay)) +  
    geom_point(alpha = 1/10)
```



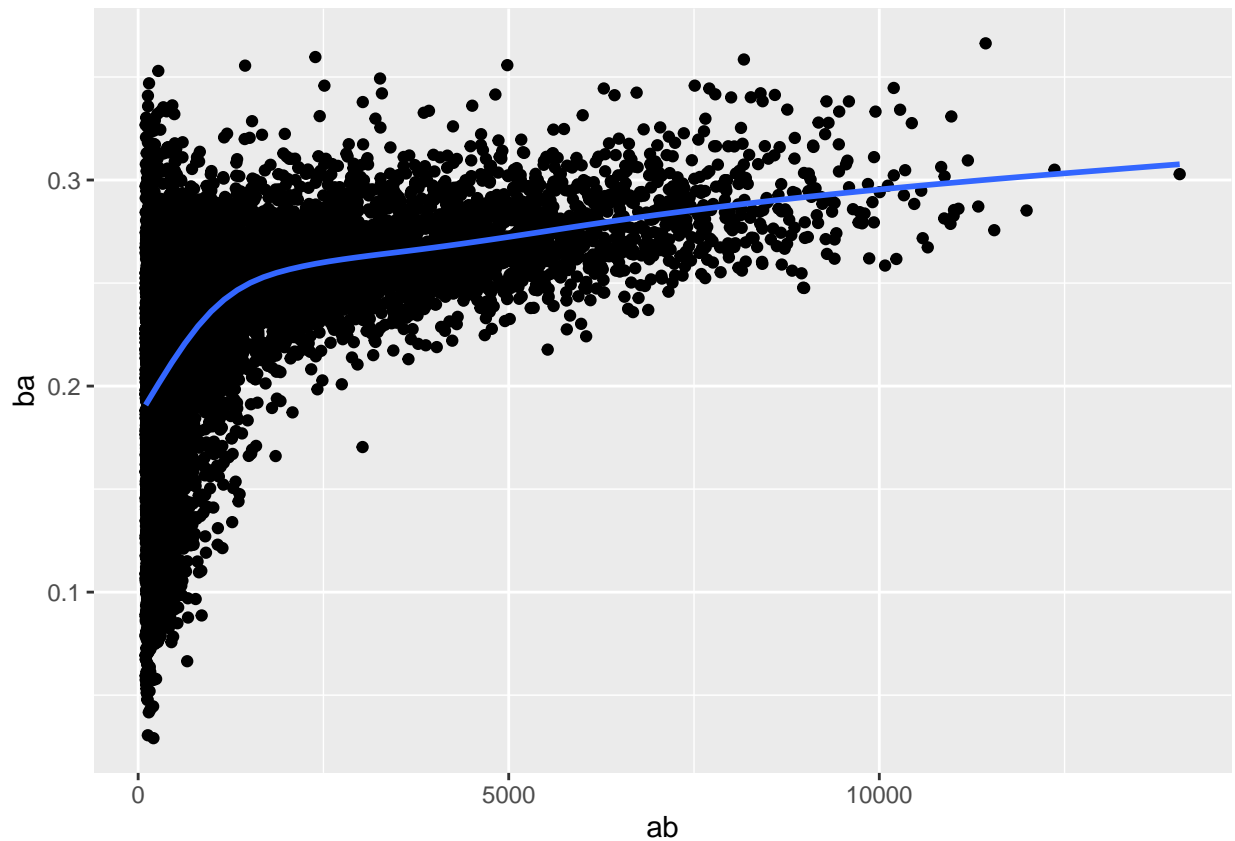
```
# Convert to a tibble so it prints nicely
batting <- as_tibble(Lahman::Batting)

batters <- batting %>%
  group_by(playerID) %>%
  summarise(
    ba = sum(H, na.rm = TRUE) / sum(AB, na.rm = TRUE),
    ab = sum(AB, na.rm = TRUE)
  )

batters %>%
  filter(ab > 100) %>%
  ggplot(mapping = aes(x = ab, y = ba)) +
  geom_point() +
  geom_smooth(se = FALSE)

## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```





```
#> `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```

```
#sort on desc(ba)
batters %>%
  arrange(desc(ba))
```

```
## # A tibble: 19,428 x 3
##   playerID    ba    ab
##   <chr>      <dbl> <int>
## 1 abramge01     1     1
## 2 alberan01     1     1
## 3 allarko01     1     1
## 4 banisje01     1     1
## 5 bartocl01     1     1
## 6 bassdo01      1     1
## 7 birasst01     1     2
## 8 bruneju01     1     1
## 9 burnscb01     1     1
## 10 cammaer01    1     1
## # ... with 19,418 more rows
```

```
#Useful summary functions
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(
    avg_delay1 = mean(arr_delay),
    avg_delay2 = mean(arr_delay[arr_delay > 0]) # the average positive delay
```

```

)

## # A tibble: 365 x 5
## # Groups:   year, month [12]
##   year month   day avg_delay1 avg_delay2
##   <int> <int> <int>     <dbl>     <dbl>
## 1  2013     1     1      12.7      32.5
## 2  2013     1     2      12.7      32.0
## 3  2013     1     3       5.73     27.7
## 4  2013     1     4      -1.93     28.3
## 5  2013     1     5      -1.53     22.6
## 6  2013     1     6       4.24     24.4
## 7  2013     1     7      -4.95     27.8
## 8  2013     1     8      -3.23     20.8
## 9  2013     1     9      -0.264    25.6
## 10 2013     1    10      -5.90     27.3
## # ... with 355 more rows

# Why is distance to some destinations more variable than to others?
not_cancelled %>%
  group_by(dest) %>%
  summarise(distance_sd = sd(distance)) %>%
  arrange(desc(distance_sd))

## # A tibble: 104 x 2
##   dest distance_sd
##   <chr>         <dbl>
## 1 EGE          10.5
## 2 SAN          10.4
## 3 SFO          10.2
## 4 HNL          10.0
## 5 SEA           9.98
## 6 LAS           9.91
## 7 PDX           9.87
## 8 PHX           9.86
## 9 LAX           9.66
## 10 IND          9.46
## # ... with 94 more rows

# When do the first and last flights leave each day?
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(
    first = min(dep_time),
    last = max(dep_time)
  )

## # A tibble: 365 x 5
## # Groups:   year, month [12]
##   year month   day first last
##   <int> <int> <int> <dbl> <dbl>
## 1  2013     1     1   517 2356
## 2  2013     1     2    42 2354
## 3  2013     1     3    32 2349
## 4  2013     1     4    25 2358

```

```
## 5 2013 1 5 14 2357
## 6 2013 1 6 16 2355
## 7 2013 1 7 49 2359
## 8 2013 1 8 454 2351
## 9 2013 1 9 2 2252
## 10 2013 1 10 3 2320
## # ... with 355 more rows
```

*#Find first and last departure for each day.*

```
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(
    first_dep = first(dep_time),
    last_dep = last(dep_time)
  )
```

```
## # A tibble: 365 x 5
## # Groups:   year, month [12]
##   year month   day first_dep last_dep
##   <int> <int> <int>     <int>     <int>
## 1 2013     1     1       517       2356
## 2 2013     1     2        42       2354
## 3 2013     1     3        32       2349
## 4 2013     1     4        25       2358
## 5 2013     1     5        14       2357
## 6 2013     1     6        16       2355
## 7 2013     1     7        49       2359
## 8 2013     1     8       454       2351
## 9 2013     1     9         2       2252
## 10 2013     1    10         3       2320
## # ... with 355 more rows
```

*# Filtering gives you all variables, with each observation in a separate row.*

```
not_cancelled %>%
  group_by(year, month, day) %>%
  mutate(r = min_rank(desc(dep_time))) %>%
  filter(r %in% range(r))
```

```
## # A tibble: 770 x 20
## # Groups:   year, month, day [365]
##   year month   day dep_time sched_dep_time dep_delay arr_time
##   <int> <int> <int>     <int>         <int>      <dbl>     <int>
## 1 2013     1     1       517           515         2       830
## 2 2013     1     1      2356          2359        -3       425
## 3 2013     1     2        42          2359        43       518
## 4 2013     1     2      2354          2359        -5       413
## 5 2013     1     3        32          2359        33       504
## 6 2013     1     3      2349          2359       -10       434
## 7 2013     1     4        25          2359        26       505
## 8 2013     1     4      2358          2359        -1       429
## 9 2013     1     4      2358          2359        -1       436
## 10 2013     1     5        14          2359        15       503
## # ... with 760 more rows, and 13 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 <dtm>, r <int>
```

```
# Which destinations have the most carriers?
not_cancelled %>%
  group_by(dest) %>%
  summarise(carriers = n_distinct(carrier)) %>%
  arrange(desc(carriers))
```

```
## # A tibble: 104 x 2
##   dest carriers
##   <chr>     <int>
## 1 ATL         7
## 2 BOS         7
## 3 CLT         7
## 4 ORD         7
## 5 TPA         7
## 6 AUS         6
## 7 DCA         6
## 8 DTW         6
## 9 IAD         6
## 10 MSP        6
## # ... with 94 more rows
```

```
#Dplyr provides a simple helper count
not_cancelled %>%
  count(dest)
```

```
## # A tibble: 104 x 2
##   dest      n
##   <chr> <int>
## 1 ABQ    254
## 2 ACK    264
## 3 ALB    418
## 4 ANC      8
## 5 ATL  16837
## 6 AUS   2411
## 7 AVL    261
## 8 BDL    412
## 9 BGR    358
## 10 BHM    269
## # ... with 94 more rows
```

```
# Using weight variable.
not_cancelled %>%
  count(tailnum, wt = distance)
```

```
## # A tibble: 4,037 x 2
##   tailnum      n
##   <chr>   <dbl>
## 1 D942DN   3418
## 2 NOEGMQ  239143
## 3 N10156  109664
## 4 N102UW   25722
## 5 N103US   24619
## 6 N104UW   24616
## 7 N10575  139903
## 8 N105UW   23618
## 9 N107US   21677
```

```
## 10 N108UW 32070
## # ... with 4,027 more rows

# How many flights left before 5am? (these usually indicate delayed
# flights from the previous day)
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(n_early = sum(dep_time < 500))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day n_early
##   <int> <int> <int>   <int>
## 1  2013     1     1         0
## 2  2013     1     2         3
## 3  2013     1     3         4
## 4  2013     1     4         3
## 5  2013     1     5         3
## 6  2013     1     6         2
## 7  2013     1     7         2
## 8  2013     1     8         1
## 9  2013     1     9         3
## 10 2013     1    10         3
## # ... with 355 more rows
```

```
# What proportion of flights are delayed by more than an hour?
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(hour_perc = mean(arr_delay > 60))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day hour_perc
##   <int> <int> <int>   <dbl>
## 1  2013     1     1  0.0722
## 2  2013     1     2  0.0851
## 3  2013     1     3  0.0567
## 4  2013     1     4  0.0396
## 5  2013     1     5  0.0349
## 6  2013     1     6  0.0470
## 7  2013     1     7  0.0333
## 8  2013     1     8  0.0213
## 9  2013     1     9  0.0202
## 10 2013     1    10  0.0183
## # ... with 355 more rows
```

```
# Grouping by multiple variables
daily <- group_by(flights, year, month, day)
(per_day <- summarise(daily, flights = n()))
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##   year month   day flights
##   <int> <int> <int>   <int>
## 1  2013     1     1     842
## 2  2013     1     2     943
```

```
## 3 2013 1 3 914
## 4 2013 1 4 915
## 5 2013 1 5 720
## 6 2013 1 6 832
## 7 2013 1 7 933
## 8 2013 1 8 899
## 9 2013 1 9 902
## 10 2013 1 10 932
## # ... with 355 more rows
```

```
(per_month <- summarise(per_day, flights = sum(flights)))
```

```
## # A tibble: 12 x 3
## # Groups:   year [1]
##   year month flights
##   <int> <int>   <int>
## 1 2013     1  27004
## 2 2013     2  24951
## 3 2013     3  28834
## 4 2013     4  28330
## 5 2013     5  28796
## 6 2013     6  28243
## 7 2013     7  29425
## 8 2013     8  29327
## 9 2013     9  27574
## 10 2013    10  28889
## 11 2013    11  27268
## 12 2013    12  28135
```

```
(per_year <- summarise(per_month, flights = sum(flights)))
```

```
## # A tibble: 1 x 2
##   year flights
##   <int>   <int>
## 1 2013  336776
```

```
# To remove grouping use Ungrouping
```

```
daily %>%
  ungroup() %>%           # no longer grouped by date
  summarise(flights = n()) # all flights
```

```
## # A tibble: 1 x 1
##   flights
##   <int>
## 1 336776
```

```
#Grouped mutates (and filters)
#Find the worst members of each group
flights_sml %>%
  group_by(year, month, day) %>%
  filter(rank(desc(arr_delay)) < 10)
```

```
## # A tibble: 3,306 x 7
## # Groups:   year, month, day [365]
##   year month day dep_delay arr_delay distance air_time
##   <int> <int> <int>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 2013     1     1     853     851     184     41
```

```
## 2 2013 1 1 290 338 1134 213
## 3 2013 1 1 260 263 266 46
## 4 2013 1 1 157 174 213 60
## 5 2013 1 1 216 222 708 121
## 6 2013 1 1 255 250 589 115
## 7 2013 1 1 285 246 1085 146
## 8 2013 1 1 192 191 199 44
## 9 2013 1 1 379 456 1092 222
## 10 2013 1 2 224 207 550 94
## # ... with 3,296 more rows
```

```
# To Find all groups bigger than a threshold
```

```
popular_dests <- flights %>%
```

```
  group_by(dest) %>%
```

```
  filter(n() > 365)
```

```
popular_dests
```

```
## # A tibble: 332,577 x 19
```

```
## # Groups:   dest [77]
```

```
##   year month   day dep_time sched_dep_time dep_delay arr_time
```

```
##   <int> <int> <int>   <int>         <int>      <dbl>   <int>
```

```
## 1 2013     1     1     517           515         2     830
```

```
## 2 2013     1     1     533           529         4     850
```

```
## 3 2013     1     1     542           540         2     923
```

```
## 4 2013     1     1     544           545        -1    1004
```

```
## 5 2013     1     1     554           600        -6     812
```

```
## 6 2013     1     1     554           558        -4     740
```

```
## 7 2013     1     1     555           600        -5     913
```

```
## 8 2013     1     1     557           600        -3     709
```

```
## 9 2013     1     1     557           600        -3     838
```

```
## 10 2013     1     1     558           600        -2     753
```

```
## # ... with 332,567 more rows, and 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 <dtm>
```

```
# Standardise to compute per group metrics
```

```
popular_dests %>%
```

```
  filter(arr_delay > 0) %>%
```

```
  mutate(prop_delay = arr_delay / sum(arr_delay)) %>%
```

```
  select(year:day, dest, arr_delay, prop_delay)
```

```
## # A tibble: 131,106 x 6
```

```
## # Groups:   dest [77]
```

```
##   year month   day dest   arr_delay prop_delay
```

```
##   <int> <int> <int> <chr>      <dbl>      <dbl>
```

```
## 1 2013     1     1 IAH         11 0.000111
```

```
## 2 2013     1     1 IAH         20 0.000201
```

```
## 3 2013     1     1 MIA         33 0.000235
```

```
## 4 2013     1     1 ORD         12 0.0000424
```

```
## 5 2013     1     1 FLL         19 0.0000938
```

```
## 6 2013     1     1 ORD          8 0.0000283
```

```
## 7 2013     1     1 LAX          7 0.0000344
```

```
## 8 2013     1     1 DFW         31 0.000282
```

```
## 9 2013     1     1 ATL         12 0.0000400
```

```
## 10 2013     1     1 DTW         16 0.000116
```

```
## # ... with 131,096 more rows
```