

day code 8

Anil Kumar Yadav

10/7/2019

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

```
##      speed      dist
##  Min.   : 4.0    Min.   :  2.00
## 1st Qu.:12.0    1st Qu.: 26.00
##  Median :15.0    Median : 36.00
##   Mean  :15.4    Mean   : 42.98
## 3rd Qu.:19.0    3rd Qu.: 56.00
##   Max.  :25.0    Max.    :120.00
```

Including Plots

You can also embed plots, for example:



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

Loading required libraries

```
library(tidyverse)
```

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

```
## v ggplot2 3.2.1    v purrr  0.3.2
## v tibble  2.1.3    v dplyr  0.8.3
## v tidyr   1.0.0    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.4.0
```

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

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

```
library(ggplot2)
library(data.table)
```

```
##
## Attaching package: 'data.table'
```

```
## The following objects are masked from 'package:dplyr':
##
##   between, first, last
```

```
## The following object is masked from 'package:purrr':
##
##   transpose
```

```
library(nycflights13)
library(Lahman)
```

Understanding the data (data validation)

```
data("flights")
head(flights)
```

```
## # A tibble: 6 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
## # ... with 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>
```

```
str(flights)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':   336776 obs. of  19 variables:
##  $ year      : int  2013 2013 2013 2013 2013 2013 2013 2013 2013 2013 ...
##  $ month     : int   1  1  1  1  1  1  1  1  1  1 ...
##  $ day       : int   1  1  1  1  1  1  1  1  1  1 ...
##  $ dep_time  : int  517 533 542 544 554 554 555 557 557 558 ...
##  $ sched_dep_time: int  515 529 540 545 600 558 600 600 600 600 ...
##  $ dep_delay : num   2  4  2 -1 -6 -4 -5 -3 -3 -2 ...
##  $ arr_time  : int  830 850 923 1004 812 740 913 709 838 753 ...
##  $ sched_arr_time: int  819 830 850 1022 837 728 854 723 846 745 ...
##  $ arr_delay  : num  11 20 33 -18 -25 12 19 -14 -8 8 ...
##  $ carrier   : chr  "UA" "UA" "AA" "B6" ...
##  $ flight    : int 1545 1714 1141 725 461 1696 507 5708 79 301 ...
##  $ tailnum   : chr  "N14228" "N24211" "N619AA" "N804JB" ...
##  $ origin    : chr  "EWR" "LGA" "JFK" "JFK" ...
##  $ dest      : chr  "IAH" "IAH" "MIA" "BQN" ...
##  $ air_time  : num  227 227 160 183 116 150 158 53 140 138 ...
##  $ distance  : num 1400 1416 1089 1576 762 ...
##  $ hour      : num   5  5  5  5  6  5  6  6  6 ...
##  $ minute    : num  15 29 40 45  0 58  0  0  0 ...
##  $ time_hour : POSIXct, format: "2013-01-01 05:00:00" "2013-01-01 05:00:00" ...
```

```
summary(flights)
```

```
##      year      month      day      dep_time
## Min.   :2013   Min.   : 1.000   Min.   : 1.00   Min.   :    1
## 1st Qu.:2013   1st Qu.: 4.000   1st Qu.: 8.00   1st Qu.: 907
## Median :2013   Median : 7.000   Median :16.00   Median :1401
## Mean   :2013   Mean   : 6.549   Mean   :15.71   Mean   :1349
## 3rd Qu.:2013   3rd Qu.:10.000   3rd Qu.:23.00   3rd Qu.:1744
## Max.   :2013   Max.   :12.000   Max.   :31.00   Max.   :2400
##                                     NA's   :8255
## sched_dep_time  dep_delay      arr_time  sched_arr_time
## Min.   : 106   Min.   : -43.00   Min.   :    1   Min.   :    1
## 1st Qu.: 906   1st Qu.:  -5.00   1st Qu.:1104   1st Qu.:1124
## Median :1359   Median :  -2.00   Median :1535   Median :1556
## Mean   :1344   Mean   : 12.64   Mean   :1502   Mean   :1536
## 3rd Qu.:1729   3rd Qu.: 11.00   3rd Qu.:1940   3rd Qu.:1945
## Max.   :2359   Max.   :1301.00   Max.   :2400   Max.   :2359
##                                     NA's   :8255   NA's   :8713
##      arr_delay      carrier      flight      tailnum
## Min.   : -86.000   Length:336776   Min.   :    1   Length:336776
## 1st Qu.: -17.000   Class :character 1st Qu.: 553   Class :character
## Median :  -5.000   Mode  :character Median :1496   Mode  :character
## Mean    :   6.895                                     Mean   :1972
## 3rd Qu.: 14.000                                     3rd Qu.:3465
## Max.    :1272.000                                    Max.   :8500
## NA's    :9430
##      origin      dest      air_time      distance
## Length:336776   Length:336776   Min.   : 20.0   Min.   : 17
## Class :character Class :character 1st Qu.: 82.0   1st Qu.: 502
## Mode  :character Mode  :character Median :129.0   Median : 872
##                                     Mean   :150.7   Mean   :1040
##                                     3rd Qu.:192.0   3rd Qu.:1389
##                                     Max.   :695.0   Max.   :4983
##                                     NA's   :9430
##      hour      minute      time_hour
## Min.   : 1.00   Min.   : 0.00   Min.   :2013-01-01 05:00:00
## 1st Qu.: 9.00   1st Qu.: 8.00   1st Qu.:2013-04-04 13:00:00
## Median :13.00   Median :29.00   Median :2013-07-03 10:00:00
## Mean   :13.18   Mean   :26.23   Mean   :2013-07-03 05:22:54
## 3rd Qu.:17.00   3rd Qu.:44.00   3rd Qu.:2013-10-01 07:00:00
## Max.   :23.00   Max.   :59.00   Max.   :2013-12-31 23:00:00
##
```

Use of filter function

```
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>
```

We can pick the right observations by using the function, `filter()`

```
jan1 <- filter(flights, month==1, day==1)
```

As we know, all the data is stored in the variable, “jan1”

```
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>
```

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

The `%in%` uses the values from the month of 11 & 12 , i.e nov and dec. So you can use either operator (`%in%`) or use the logical operator in the code

```
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>
```

Both the lines of code same output. But we can use either of it based on our convinence

Using arrange() function

```
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>
```

```
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>
```

The function, desc() gives the values of the mentioned variable in the decreasing fashion

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

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

As you can the NA values are always stored in the bottom

```
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() is used select the exclusive data set

```
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>
```

We use the rename() function because it not only picks the mentioned list but also the variables that are not mentioned specifically

```
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
```

With the mutate() we had added two new variables i.e., "gain", "speed"


```
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
```

Use the function `transmute()` to keep up the new variables along with the given data

```
summarise(flights, delay=mean(dep_delay, na.rm = TRUE))
```

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

It summarised the whole data

```
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
```

here we have group the data first and later found out the mean values of the column and summarised it to given value respectively

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

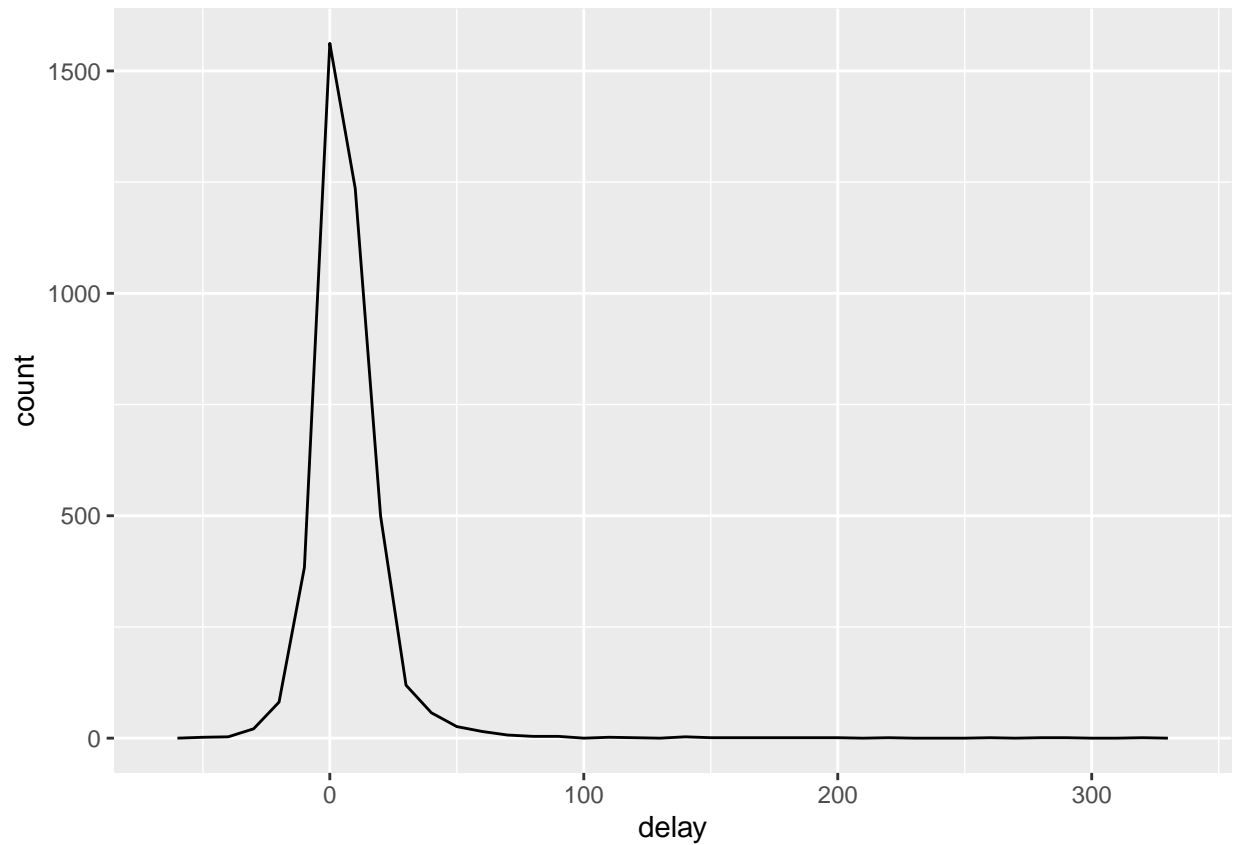
```
not_cancelled %>%
  group_by(year, month, day) %>%
  summarise(mean = mean(dep_delay))
```

```
## # 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
```

```
delays <- not_cancelled %>%
  group_by(tailnum) %>%
  summarise(delay= mean(arr_delay))
```

plotting

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



```
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 %>%
  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
```