

Instructions

Exercises: 1-3 (5.2.4 Exercises); 2-4 (5.3.1 Exercises); 2,4 (5.4.1 Exercises); 1-4 (5.5.2 Exercises)

Submission: Submit an electronic document on Gradescope. Must be submitted as a PDF file generated in RStudio. All assigned problems are chosen according to the textbook *R for Data Science*. You do not need R code to answer every question. If you answer without using R code, delete the code chunk. If the question requires R code, make sure you display R code. If the question requires a figure, make sure you display a figure. A lot of the questions can be answered in written response, but require R code and/or figures for understanding and explaining.

Chapter 5 (5.2.4 Exercises)

Exercise 1

```
#1
filter(flights, 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
## # ... with 10,190 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#2
filter(flights, dest == "IAH" | dest == "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
## # ... with 9,303 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#3
filter(flights, carrier == "UA" | carrier == "AA" | carrier == "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
## # ... with 139,494 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#4
filter(flights, month >=7 & month <=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
## # ... with 86,316 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#5
filter(flights, arr_delay > 120, dep_delay <=0 )
```

```
## # A tibble: 29 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    27    1419           1420          -1     1754           1550
## 2  2013    10     7    1350           1350           0     1736           1526
## 3  2013    10     7    1357           1359          -2     1858           1654
## 4  2013    10    16     657           700          -3     1258           1056
## 5  2013    11     1     658           700          -2     1329           1015
## 6  2013     3    18    1844           1847          -3         39           2219
## 7  2013     4    17    1635           1640          -5     2049           1845
## 8  2013     4    18     558           600          -2     1149            850
## 9  2013     4    18     655           700          -5     1213            950
## 10 2013     5    22    1827           1830          -3     2217           2010
## # ... with 19 more rows, and 11 more variables: arr_delay <dbl>, carrier <chr>,
## #   flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
## #   distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#6
filter(flights,)
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2     830           819
## 2  2013     1     1     533           529           4     850           830
## 3  2013     1     1     542           540           2     923           850
## 4  2013     1     1     544           545          -1    1004          1022
## 5  2013     1     1     554           600          -6     812           837
## 6  2013     1     1     554           558          -4     740           728
## 7  2013     1     1     555           600          -5     913           854
## 8  2013     1     1     557           600          -3     709           723
## 9  2013     1     1     557           600          -3     838           846
## 10 2013     1     1     558           600          -2     753           745
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#7
filter(flights, dep_time <=2400, dep_time <= 600 )
```

```
## # A tibble: 9,344 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2     830           819
## 2  2013     1     1     533           529           4     850           830
## 3  2013     1     1     542           540           2     923           850
## 4  2013     1     1     544           545          -1    1004          1022
## 5  2013     1     1     554           600          -6     812           837
## 6  2013     1     1     554           558          -4     740           728
## 7  2013     1     1     555           600          -5     913           854
## 8  2013     1     1     557           600          -3     709           723
## 9  2013     1     1     557           600          -3     838           846
## 10 2013     1     1     558           600          -2     753           745
## # ... with 9,334 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Exercise 2

?between is a shortcut that is used for >= and <=. This could have been helpful for #7 for example in the last problem. It could have been used for the time between 12am - 6am.

Exercise 3

```
#Missing depart times
filter(flights, is.na(dep_time))
```

```
## # A tibble: 8,255 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     NA           1630           NA       NA           1815
## 2  2013     1     1     NA           1935           NA       NA           2240
## 3  2013     1     1     NA           1500           NA       NA           1825
## 4  2013     1     1     NA            600           NA       NA            901
## 5  2013     1     2     NA           1540           NA       NA           1747
## 6  2013     1     2     NA           1620           NA       NA           1746
## 7  2013     1     2     NA           1355           NA       NA           1459
## 8  2013     1     2     NA           1420           NA       NA           1644
## 9  2013     1     2     NA           1321           NA       NA           1536
## 10 2013     1     2     NA           1545           NA       NA           1910
## # ... with 8,245 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#8255
filter(flights, is.na(dep_time))
```

```
## # A tibble: 8,255 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     NA           1630           NA       NA           1815
## 2  2013     1     1     NA           1935           NA       NA           2240
## 3  2013     1     1     NA           1500           NA       NA           1825
## 4  2013     1     1     NA            600           NA       NA            901
## 5  2013     1     2     NA           1540           NA       NA           1747
## 6  2013     1     2     NA           1620           NA       NA           1746
## 7  2013     1     2     NA           1355           NA       NA           1459
## 8  2013     1     2     NA           1420           NA       NA           1644
## 9  2013     1     2     NA           1321           NA       NA           1536
## 10 2013     1     2     NA           1545           NA       NA           1910
## # ... with 8,245 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Chapter 5 (5.3.1 Exercises)

Exercise 2

```
#Find the most delayed flights
arrange(flights, desc(dep_delay))
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>       <dbl>   <int>         <int>
## 1  2013     1     9     641           900       1301    1242         1530
## 2  2013     6    15    1432          1935       1137    1607         2120
## 3  2013     1    10    1121          1635       1126    1239         1810
## 4  2013     9    20    1139          1845       1014    1457         2210
## 5  2013     7    22     845          1600       1005    1044         1815
## 6  2013     4    10    1100          1900        960    1342         2211
## 7  2013     3    17    2321           810        911     135         1020
## 8  2013     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
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#Find the earliest delayed flights
arrange(flights, dep_delay)
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>       <dbl>   <int>         <int>
## 1  2013    12     7    2040          2123        -43     40         2352
## 2  2013     2     3    2022          2055        -33    2240         2338
## 3  2013    11    10    1408          1440        -32    1549         1559
## 4  2013     1    11    1900          1930        -30    2233         2243
## 5  2013     1    29    1703          1730        -27    1947         1957
## 6  2013     8     9     729           755        -26    1002          955
## 7  2013    10    23    1907          1932        -25    2143         2143
## 8  2013     3    30    2030          2055        -25    2213         2250
## 9  2013     3     2    1431          1455        -24    1601         1631
## 10 2013     5     5     934           958        -24    1225         1309
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Exercise 3

```
#Sort to find the fastest flights
arrange(flights, 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     1    16    1355           1315         40    1442           1411
## 2  2013     4    13     537           527         10     622           628
## 3  2013    12     6     922           851         31    1021           954
## 4  2013     2     3    2153           2129         24    2247           2224
## 5  2013     2     5    1303           1315        -12    1342           1411
## 6  2013     2    12    2123           2130         -7    2211           2225
## 7  2013     3     2    1450           1500        -10    1547           1608
## 8  2013     3     8    2026           1935         51    2131           2056
## 9  2013     3    18    1456           1329         87    1533           1426
## 10 2013     3    19    2226           2145         41    2305           2246
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Exercise 4

```
#travelled the farthest
arrange(flights, desc(distance))
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     857           900         -3    1516           1530
## 2  2013     1     2     909           900          9    1525           1530
## 3  2013     1     3     914           900         14    1504           1530
## 4  2013     1     4     900           900          0    1516           1530
## 5  2013     1     5     858           900         -2    1519           1530
## 6  2013     1     6    1019           900         79    1558           1530
## 7  2013     1     7    1042           900        102    1620           1530
## 8  2013     1     8     901           900          1    1504           1530
## 9  2013     1     9     641           900       1301    1242           1530
## 10 2013     1    10     859           900         -1    1449           1530
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
#travelled the shortest
arrange(flights, distance)
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     7    27      NA             106           NA        NA           245
## 2  2013     1     3    2127             2129          -2       2222         2224
## 3  2013     1     4    1240             1200          40       1333         1306
## 4  2013     1     4    1829             1615         134       1937         1721
## 5  2013     1     4    2128             2129          -1       2218         2224
## 6  2013     1     5    1155             1200          -5       1241         1306
## 7  2013     1     6    2125             2129          -4       2224         2224
## 8  2013     1     7    2124             2129          -5       2212         2224
## 9  2013     1     8    2127             2130          -3       2304         2225
## 10 2013     1     9    2126             2129          -3       2217         2224
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Chapter 5 (5.4.1 Exercises)

Exercise 2

Select() has a mini language that allows the variables to be easily based on their names. It is used to select the variables in a data set.

```
#
```

Exercise 4

The code does surprise me because time was written as "TIME" (it is in capital letters) in this code, so I wasn't sure how it received that information. In order to fix this, we can change the contains to have "TIME", ignore.case = FALSE, since the contains() default is not affected by case.

so the code would need to be: select(flights, contains("TIME", ignore.case = FALSE))

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



```
## # A tibble: 336,776 x 6
##   dep_time sched_dep_time arr_time sched_arr_time air_time time_hour
##   <int>      <int>      <int>      <int>      <dbl> <dtm>
## 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
## # ... with 336,766 more rows
```

Chapter 5 (5.5.2 Exercises)

Exercise 1

```
mutate(flights, dep_time = (dep_time %/% 100) * 60 + (dep_time %% 100), sched_dep_time =
= (sched_dep_time %/% 100) * 60 + (sched_dep_time %% 100))
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <dbl>       <dbl>      <dbl>    <int>         <int>
## 1  2013     1     1     317         315         2      830          819
## 2  2013     1     1     333         329         4      850          830
## 3  2013     1     1     342         340         2      923          850
## 4  2013     1     1     344         345        -1     1004         1022
## 5  2013     1     1     354         360        -6      812          837
## 6  2013     1     1     354         358        -4      740          728
## 7  2013     1     1     355         360        -5      913          854
## 8  2013     1     1     357         360        -3      709          723
## 9  2013     1     1     357         360        -3      838          846
## 10 2013     1     1     358         360        -2      753          745
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

Exercise 2

I thought that they would be the same because (arrival - departure) would equal the amount of air time.

I noticed that the numbers are different from each other. In order to fix this, I think that arr_time and dep_time need to be converted (using mutate) in order to change the units to minutes.

```
flights %>% mutate(flight_time = arr_time - dep_time) %>%
  select(arr_time, flight_time)
```

```
## # A tibble: 336,776 x 2
##   air_time flight_time
##   <dbl>      <int>
## 1      227         313
## 2      227         317
## 3      160         381
## 4      183         460
## 5      116         258
## 6      150         186
## 7      158         358
## 8       53         152
## 9      140         281
## 10     138         195
## # ... with 336,766 more rows
```

Exercise 3

I think that they are related because the $(\text{dep_time}) - (\text{sched_dep_time}) = \text{dep_delay}$.

```
select(flights, 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
## # ... with 336,766 more rows
```

Exercise 4

I wanted to use `min_rank` because it assigns the tied values at the same rank but it also assigns a rank that is equal to the amount of values less than the one tied value plus one.

```
head(arrange(flights, min_rank(desc(dep_delay))), 10)
```

```
## # A tibble: 10 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
## # ... with 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 <dtm>
```