# STAT167 Lab #4 - Spring 2025

Ethan Choi

2025/4/25

## Contents

## Discussion/Lab #4 instructions

This week, we will review some `ggplot2` examples for data visualization and `dplyr` examples for data transformation/manipulation.

- First, download the `rmd` file from Canvas.

- Open this `rmd` file in RStudio and click `Knit -> Knit to PDF` to render it to PDF format. You need to have `LaTex` installed on the computer to render it to PDF format. If not, you can also render it to HTML format.

- Read this `rmd` file and the rendered `pdf`/`html` file side-by-side, to see how this document was generated!

- Be sure to play with this document! Change it. Break it. Fix it. The best way to learn R Markdown (or really almost anything) is to try, fail, then find out what you did wrong.

- Read over the `dplyr` example code and the output. If you have any questions about certain functions or parameters, it is the time to ask!

- There are some exercises through out this document. Replace **INSERT_YOUR_ANSWER** with your own answers. Knit the file, and check your results.

**Please comment your R code thoroughly, and follow the R coding style guideline (https: //google.github.io/styleguide/Rguide.xml). Partial credit will be deducted for insufficient commenting or poor coding styles.**

**Lab submission guideline**

- After you completed all exercises, save your file to `FirstnameLastname-SID-lab4.rmd` and save the rendered pdf file to `FirstnameLastname-SID-lab4.pdf`. If you can not knit it to pdf, knit it to html first and then print/save it to pdf format.
- Submit **BOTH your source `rmd` file and the knitted `pdf` file** to **GradeScope**. Do NOT create a zip file.
- You can submit multiple times, you last submission will be graded.

---

# Lecture Review - data visualization with `ggplot2`

## Load the `tidyverse` package

```
# install the tidyverse package first if you have not done it yet.
#install.packages("tidyverse") # you can comment out this line after you have installed `tidyverse`

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   3.5.1     v tibble    3.2.1
## v lubridate 1.9.4     v tidyr     1.3.1
## v purrr     1.0.4
## -- 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 become error
```

## The `mpg` data set

The `mpg` data set contains fuel economy data 1999 - 2008 for 38 popular car models. https://ggplot2. tidyverse.org/reference/mpg.html

```
?mpg
## starting httpd help server ... done
glimpse(mpg) # get a glimpse of the mpg data
## Rows: 234
## Columns: 11
## $ manufacturer <chr> "audi", "audi", "audi", "audi", "audi", "audi", "audi", "~
## $ model        <chr> "a4", "a4", "a4", "a4", "a4", "a4", "a4", "a4 quattro", "~
## $ displ        <dbl> 1.8, 1.8, 2.0, 2.0, 2.8, 2.8, 3.1, 1.8, 1.8, 2.0, 2.0, 2.~
## $ year         <int> 1999, 1999, 2008, 2008, 1999, 1999, 2008, 1999, 1999, 200~
## $ cyl          <int> 4, 4, 4, 4, 6, 6, 6, 4, 4, 4, 4, 6, 6, 6, 6, 6, 6, 8, 8, ~
## $ trans        <chr> "auto(l5)", "manual(m5)", "manual(m6)", "auto(av)", "auto~
## $ drv          <chr> "f", "f", "f", "f", "f", "f", "f", "4", "4", "4", "4", "4~
## $ cty          <int> 18, 21, 20, 21, 16, 18, 18, 18, 16, 20, 19, 15, 17, 17, 1~
## $ hwy          <int> 29, 29, 31, 30, 26, 26, 27, 26, 25, 28, 27, 25, 25, 25, 2~
## $ fl           <chr> "p", "p", "p", "p", "p", "p", "p", "p", "p", "p", "p", "p~
## $ class        <chr> "compact", "compact", "compact", "compact", "compact", "c~
```

## The complete graphing template in `ggplot2`

```
ggplot(data = <DATA>) +
  <GEOM_FUNCTION>(
     mapping = aes(<MAPPINGS>),
     stat = <STAT>,          # optional
     position = <POSITION>  # optional
  ) +
  <COORDINATE_FUNCTION> +   # optional
  <FACET_FUNCTION> +         # optional
  <SCALE_FUNCTION> +         # optional
  <THEME_FUNCTION>           # optional
```
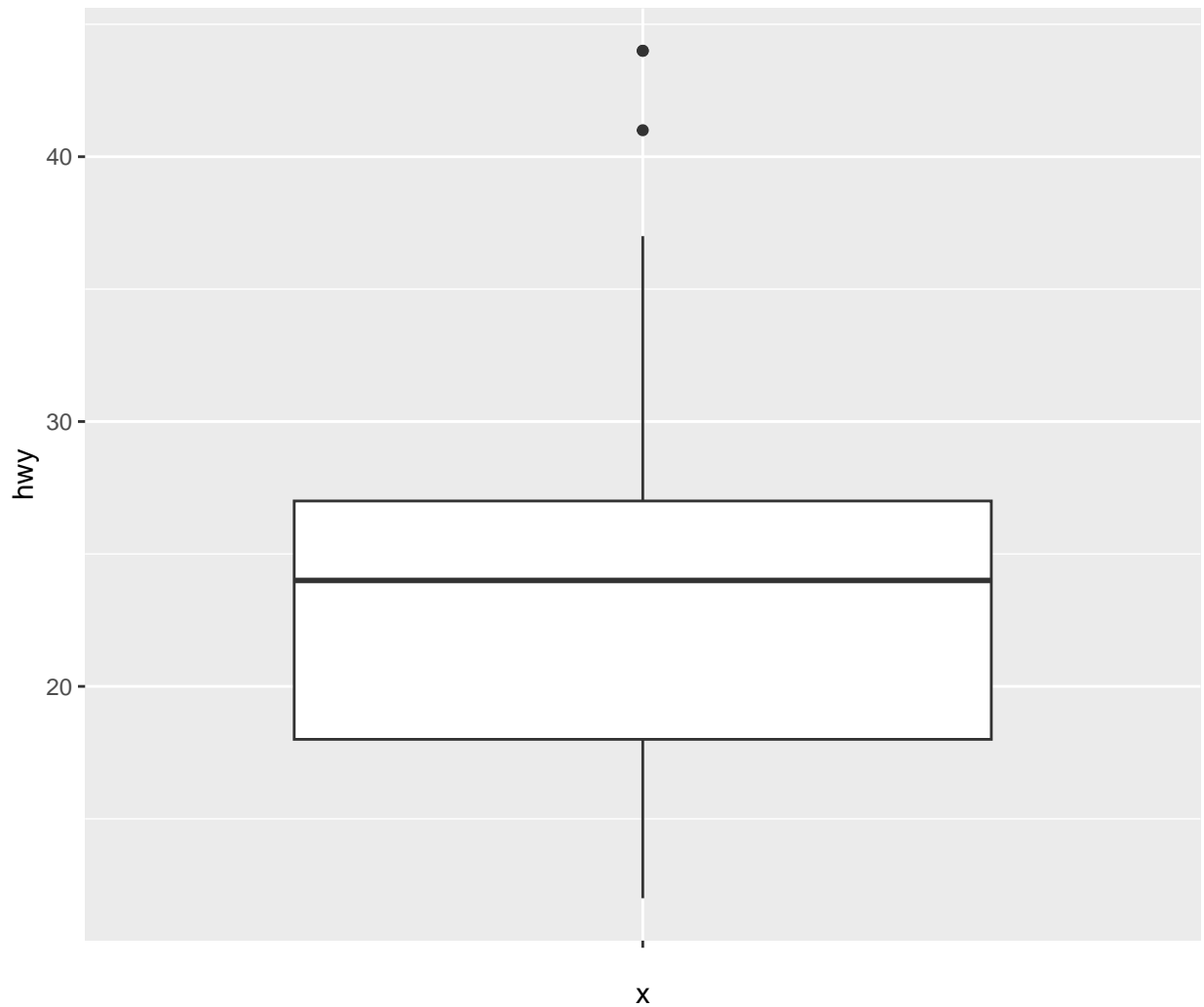
---

## Example: `geom_boxplot()` for one variable

Let's first take a look at the distribution of highway mileage.

```
ggplot(data = mpg) +
  geom_boxplot(mapping = aes(x = "", y = hwy))
```

**Exercise #1**

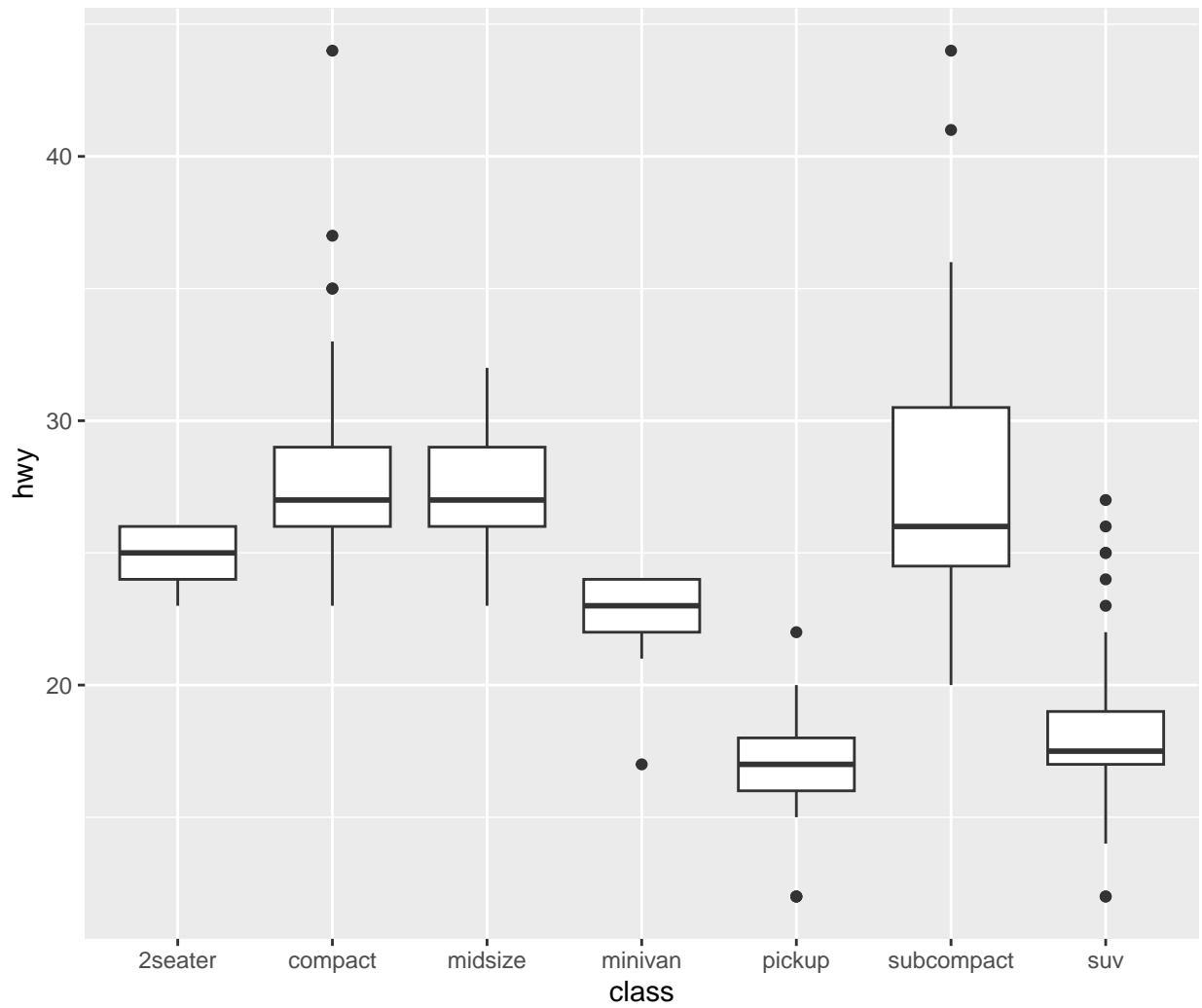If you add `facet_wrap(~class)` to the above code, what will happen?

**INSERT__YOUR__ANSWER** It will create separate boxplots of highway mileage for each car class into different individual subplots. ***

## Example: `geom_boxplot()` for two variables

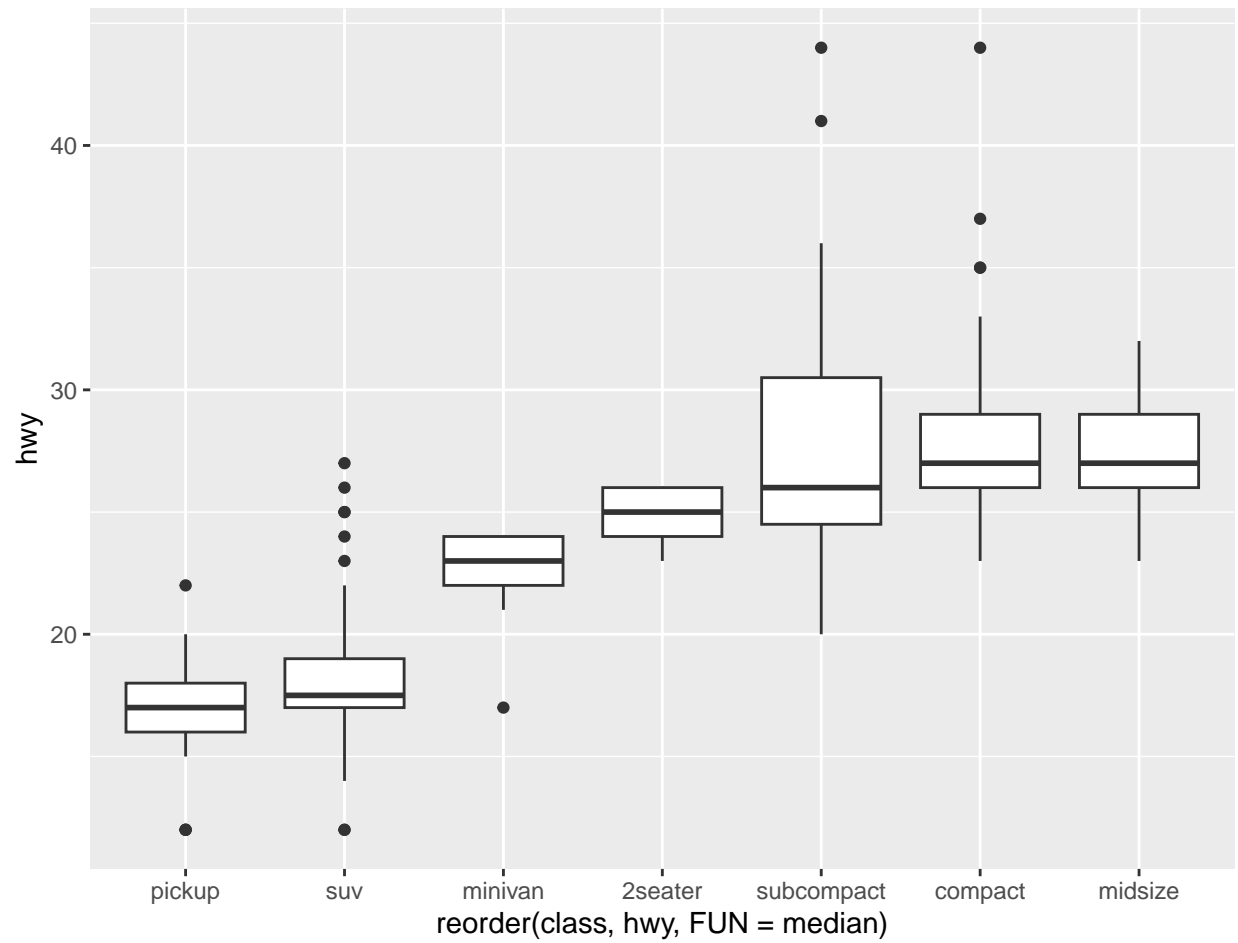How highway mileage varies across car classes?

- Plot the distribution of highway mileage by class.

```
ggplot(data = mpg) +
  geom_boxplot(mapping = aes(x = class, y = hwy))
```

To make the trend easier to see, we can reorder `class` based on the median value of `hwy`.

```
ggplot(data = mpg) +
  geom_boxplot(mapping = aes(x = reorder(class, hwy, FUN = median),
                             y = hwy))
```

When you have long variable names, `geom_boxplot()` will work better if you flip it with `coord_flip()`.

```
ggplot(data = mpg) +
  geom_boxplot(mapping = aes(x = reorder(class, hwy, FUN = median),
                             y = hwy)) +
  xlab("class") +
  coord_flip()
```

**Exercise #2**

Can you draw a violin plot to compare the distribution of highway mileage varies across car classes? Order your violins by the median values of highway mileage and flip the plot to horizontal.

**INSERT_YOUR_ANSWER**

```r
ggplot(data = mpg) +
  geom_violin(mapping = aes(x = reorder(class, hwy, FUN = median), y = hwy), fill = "lightblue") +
  coord_flip() +
  xlab("class") +
  ylab("Highway Mileage") +
  ggtitle("Violin Plot of Highway Mileage by Car Class")
```

## Violin Plot of Highway Mileage by Car Class

**Exercise #3**

There are two ways to zoom in a ggplot. Compare the following two plots. What are the differences? Which zooming option do you prefer? Explain your reasons.

**INSERT_YOUR_ANSWER** The ylim(c(#, #) method removes data points outside the range, which can distort/squish the plot and accidentally remove outliers. The coord_cartesian(ylim = c(#, #)) approach zooms into a specified range without missing any data. I prefer the second method because it keeps the full dataset intact while still focusing the view, which doesn't lead to misinterpretation of data.

```
ggplot(data = mpg) +
  geom_boxplot(mapping = aes(x = class, y = hwy)) +
  ylim(c(20, 30))
## Warning: Removed 100 rows containing non-finite outside the scale range
## (`stat_boxplot()`).
```

```
ggplot(data = mpg) +
  geom_boxplot(mapping = aes(x = class, y = hwy)) +
  coord_cartesian(ylim=c(20, 30))
```

# Lecture Review - data transformation with `dplyr`

### The `nycflights13::flights` data set

This data frame contains all 336,776 flights that departed from New York City in 2013. https://nycflights13.tidyverse.org

```
# You need to install `nycflights13` first, then you can comment out the following line
#install.packages("nycflights13")
library(nycflights13)
```

```
#?flights # full documentation
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~
```

### `filter()` - choose rows by their column values

`filter()` allows you to subset rows based on their column values.

```
# flights that departed on Jan/1/2013
filter(flights, month == 1, day == 1)
## # A tibble: 842 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013     1     1      517            515         2      830            819
## 2   2013     1     1      533            529         4      850            830
## 3   2013     1     1      542            540         2      923            850
## 4   2013     1     1      544            545        -1     1004           1022
## 5   2013     1     1      554            600        -6      812            837
## 6   2013     1     1      554            558        -4      740            728
## 7   2013     1     1      555            600        -5      913            854
## 8   2013     1     1      557            600        -3      709            723
## 9   2013     1     1      557            600        -3      838            846
```

```
## 10  2013     1     1      558          600        -2      753          745
## # i 832 more rows
## # i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>
```

```
# flights that departed in Nov or Dec
filter(flights, month == 11 | month == 12)
## # A tibble: 55,403 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013    11     1        5           2359         6      352            345
## 2   2013    11     1       35           2250       105      123           2356
## 3   2013    11     1      455            500        -5      641            651
## 4   2013    11     1      539            545        -6      856            827
## 5   2013    11     1      542            545        -3      831            855
## 6   2013    11     1      549            600       -11      912            923
## 7   2013    11     1      550            600       -10      705            659
## 8   2013    11     1      554            600        -6      659            701
## 9   2013    11     1      554            600        -6      826            827
## 10  2013    11     1      554            600        -6      749            751
## # i 55,393 more rows
## # i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>
```

```
filter(flights, month %in% c(11, 12))
## # A tibble: 55,403 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013    11     1        5           2359         6      352            345
## 2   2013    11     1       35           2250       105      123           2356
## 3   2013    11     1      455            500        -5      641            651
## 4   2013    11     1      539            545        -6      856            827
## 5   2013    11     1      542            545        -3      831            855
## 6   2013    11     1      549            600       -11      912            923
## 7   2013    11     1      550            600       -10      705            659
## 8   2013    11     1      554            600        -6      659            701
## 9   2013    11     1      554            600        -6      826            827
## 10  2013    11     1      554            600        -6      749            751
## # i 55,393 more rows
## # i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>
```

```
# flights that weren't delayed (on arrival or departure) by more than 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 sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013     1     1      517            515         2      830            819
## 2   2013     1     1      533            529         4      850            830
## 3   2013     1     1      542            540         2      923            850
```

```
##  4  2013    1    1    544           545          -1    1004           1022
##  5  2013    1    1    554           600          -6     812            837
##  6  2013    1    1    554           558          -4     740            728
##  7  2013    1    1    555           600          -5     913            854
##  8  2013    1    1    557           600          -3     709            723
##  9  2013    1    1    557           600          -3     838            846
## 10  2013    1    1    558           600          -2     753            745
## # i 316,040 more rows
## # i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>

filter(flights, arr_delay <= 120, dep_delay <= 120)
## # A tibble: 316,050 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
##  1  2013    1    1    517           515           2     830            819
##  2  2013    1    1    533           529           4     850            830
##  3  2013    1    1    542           540           2     923            850
##  4  2013    1    1    544           545          -1    1004           1022
##  5  2013    1    1    554           600          -6     812            837
##  6  2013    1    1    554           558          -4     740            728
##  7  2013    1    1    555           600          -5     913            854
##  8  2013    1    1    557           600          -3     709            723
##  9  2013    1    1    557           600          -3     838            846
## 10  2013    1    1    558           600          -2     753            745
## # i 316,040 more rows
## # i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>
```

**Exercise #4**

(a) Find the flights that departed in the summer (July, August, and September)

**INSERT_YOUR_ANSWER**

```
summer_flights <- filter(flights, month %in% c(7, 8, 9))
head(summer_flights)
```

```
## # A tibble: 6 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
## # 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>
```

(b) Find the flights that flew to Houston (IAH or HOU)

**INSERT_YOUR_ANSWER**

```
houston_flights <- filter(flights, dest %in% c("IAH", "HOU"))
head(houston_flights)
```

```
## # A tibble: 6 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
## # 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>
```