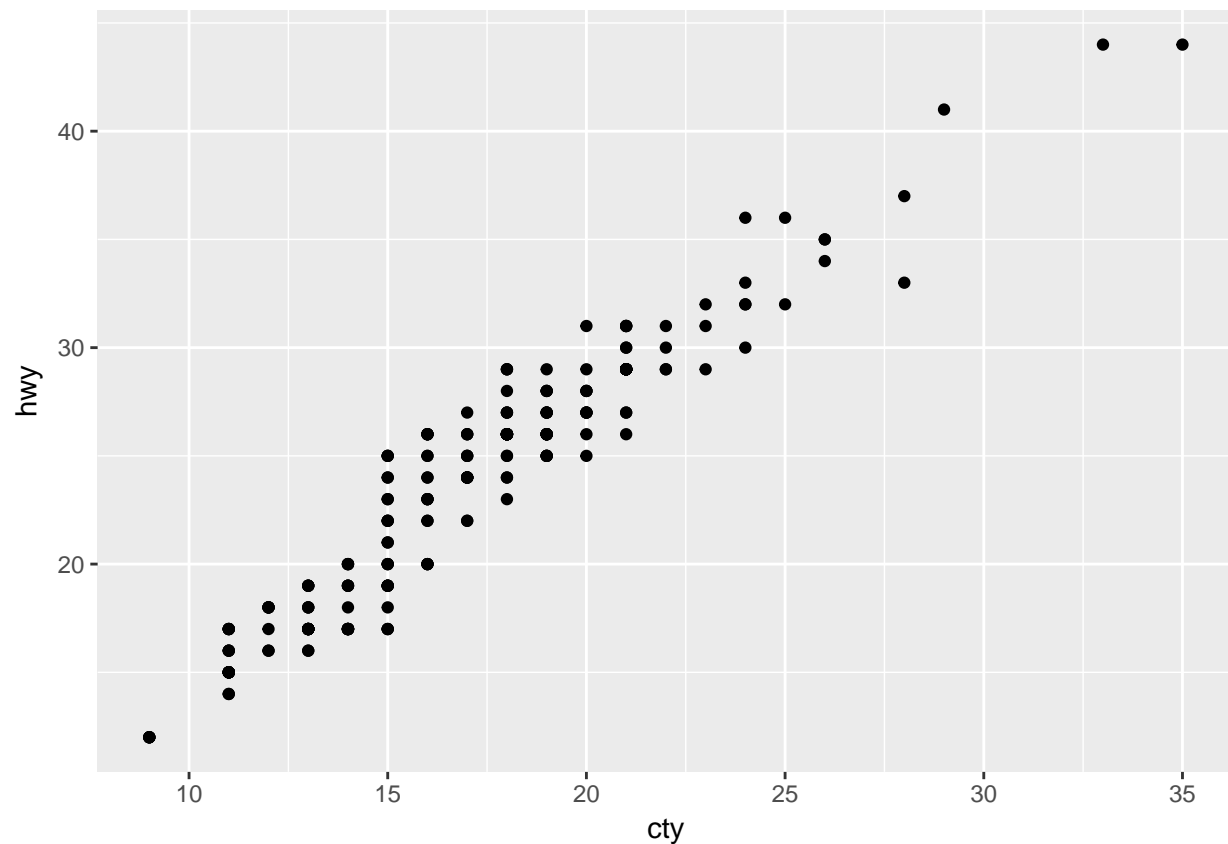# Home Work Assignment 3

*Pradeep Sahoo*

*6/21/2018*

My Github repository for my assignments can be found at this URL: My Github

**Exercises**

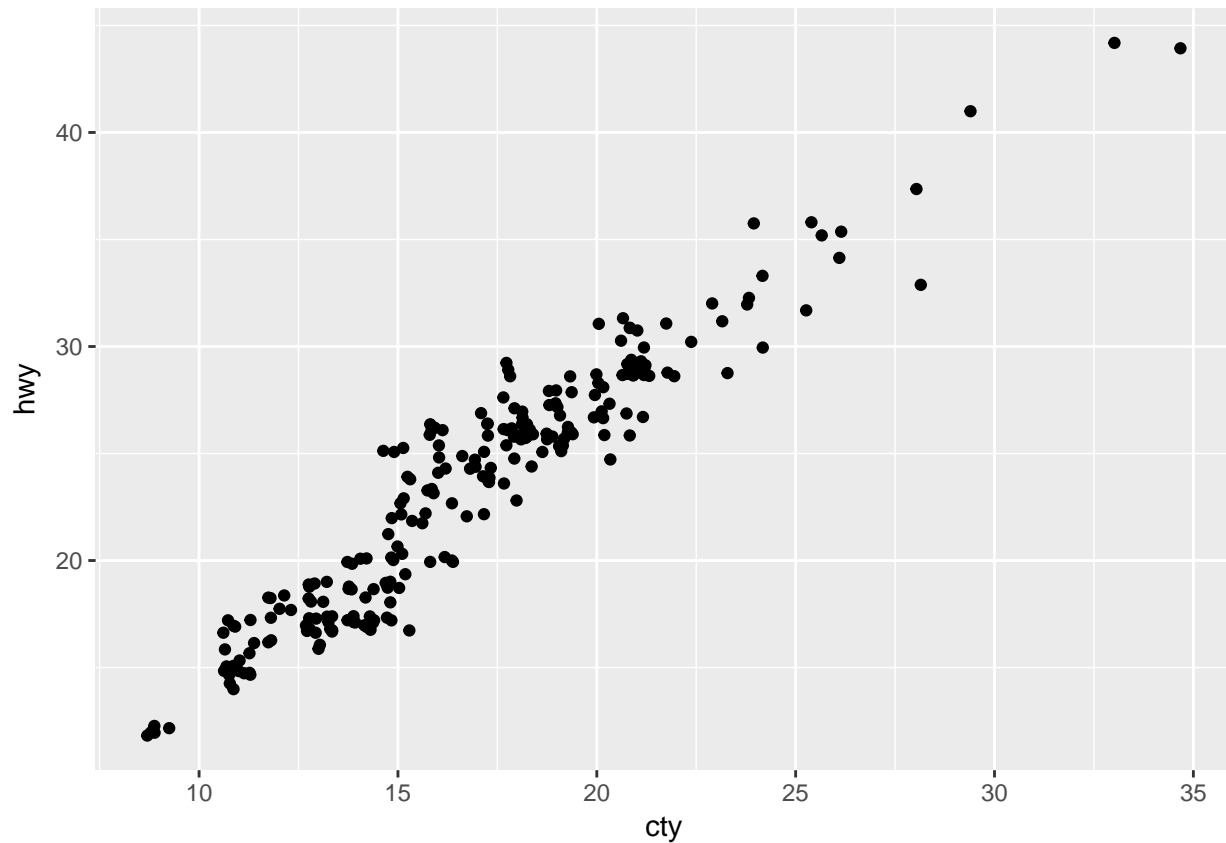## Section 3.8.1: all exercises

```
library(mdsr)
library(tidyverse)
library(nycflights13)
```

```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_point()
```



By selecting Hwy and Cty we are plotting continous variable for X and Y. To get some information out of it we should have Positions. That will help to segregate these continous variable.
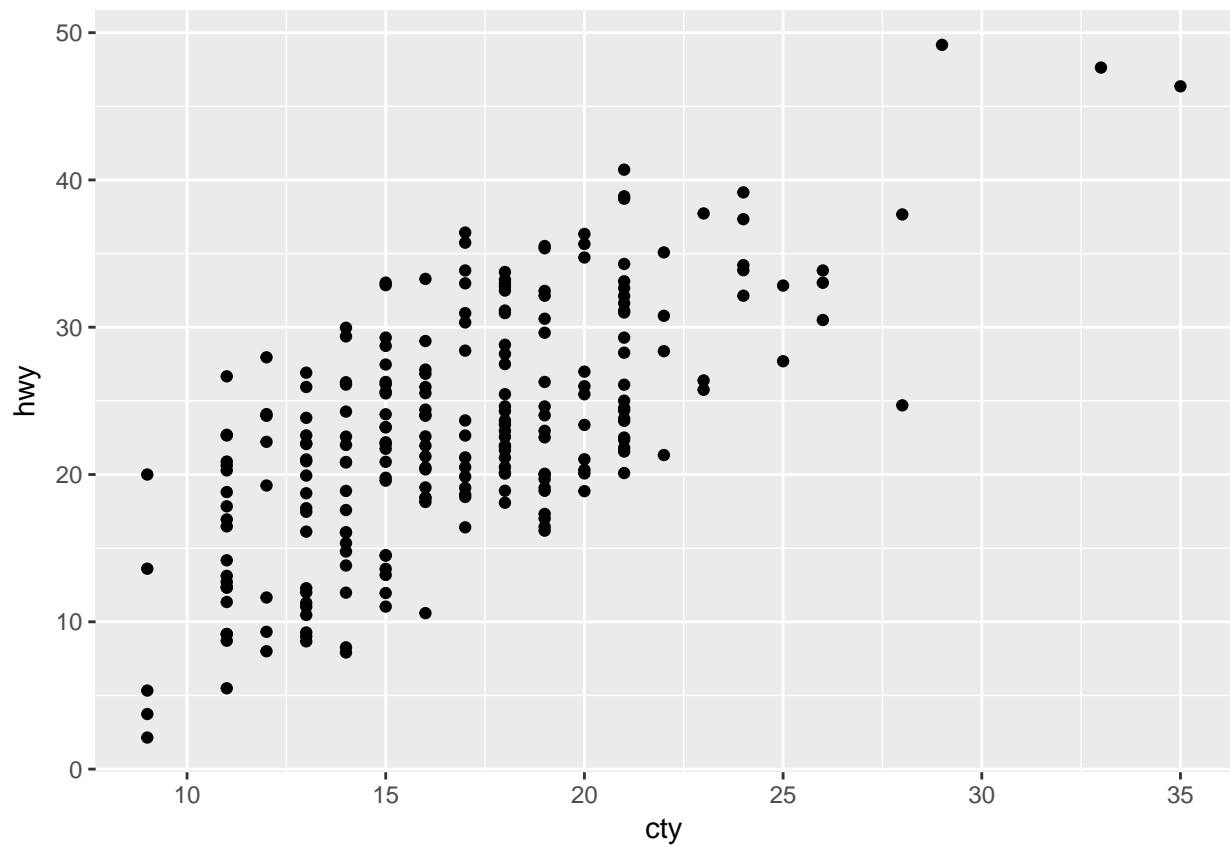
```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_point(position = "jitter")
```
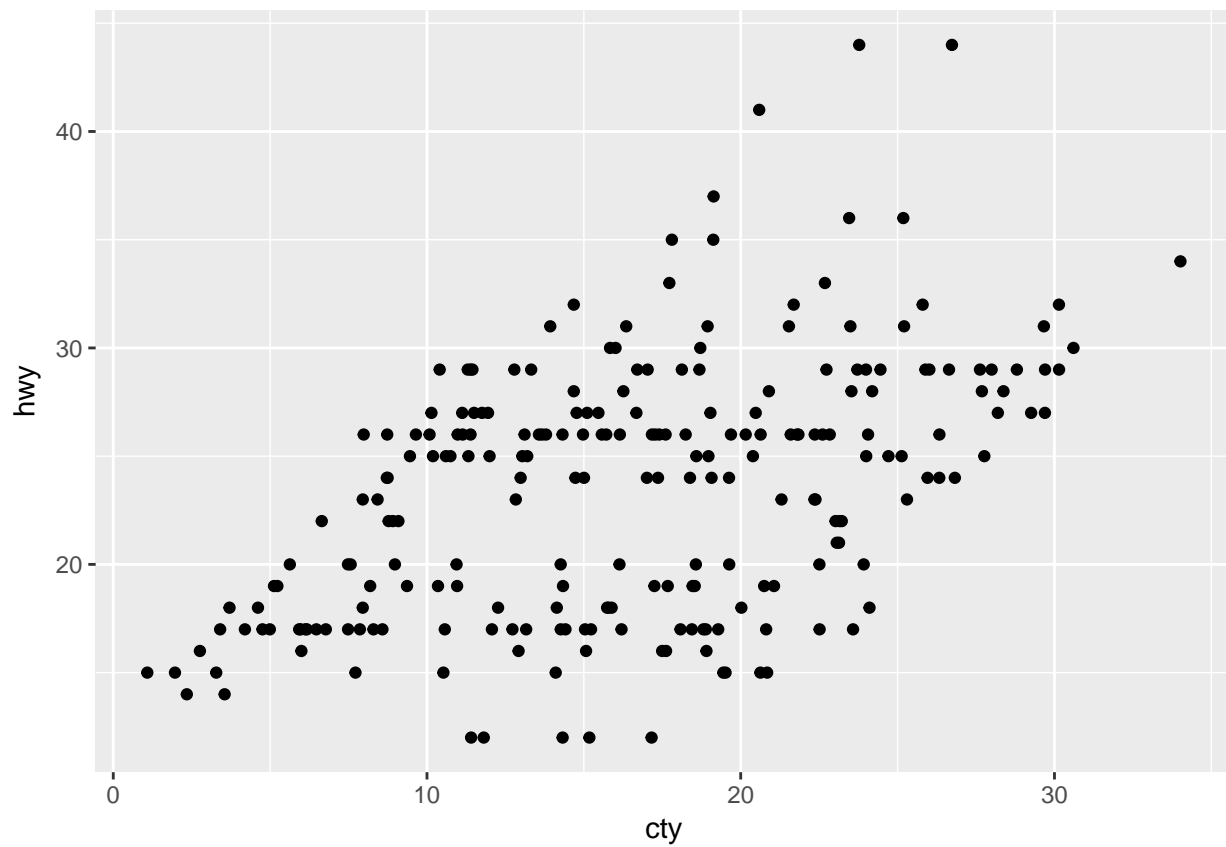
**giom_jitter function adds some amount of random variation to the location of each point. This is way to handling overplotting caused by discretness in datasets.

This has 2 parameters, width and height. Width gives a x axis variation, Height gives a Y axis variation.

```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_jitter(width = 0, height = 10)
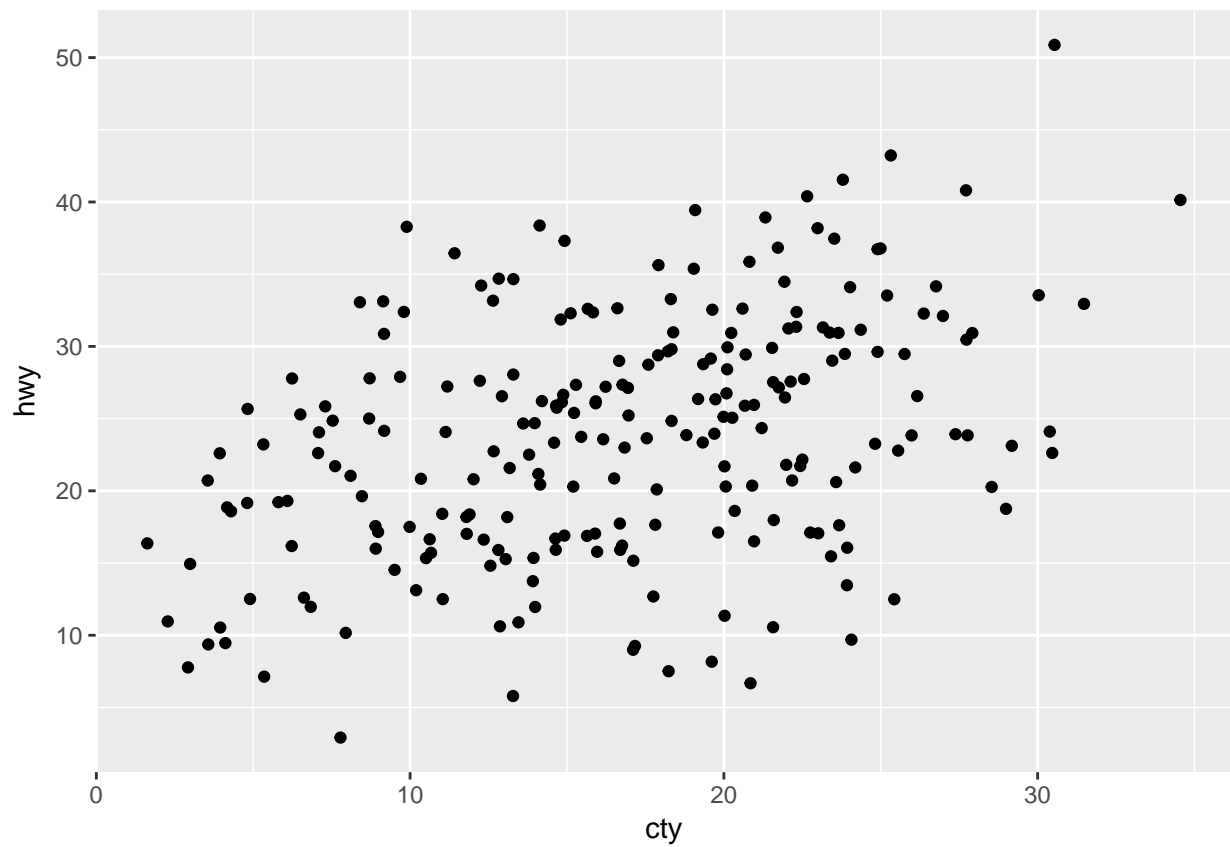```

```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_jitter(width = 10, height = 0)
```

```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_jitter(width = 10, height = 10)
```

```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_count(position = 'jitter')
```

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

```
ggplot(data = mpg, mapping = aes(x = class, y = hwy, color = drv)) +
  geom_boxplot(position = 'identity')
```

## Section 3.9.1: #2 and #4 only

**1**

```r
ggplot(mpg, aes(x = class, fill = drv)) +
  geom_bar()
```

```
ggplot(mpg, aes(x = class, fill = drv)) +
  geom_bar() + coord_polar()
```

**4**

```
ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
  geom_point() +
  geom_abline() +
  coord_fixed()
```

coord_fixed() - This function fixes the line exactly to the points where cty mileage and hwy mileage matches. It draws the line connecting those points. This gives a perspective to the occurances of the Data.

geom_abline() - This function draws a reference line. This reference line is used to give inference about the data occurances.

## Section 4.4: #1 and #2 only

## 1

Spelling Mistake.

```
my_variable <- 10
my_variable
```

```
## [1] 10
```

```
library(tidyverse)

ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy))
```

```
filter(mpg, cyl == 8)
```

```
## # A tibble: 70 x 11
##    manufacturer model      displ  year   cyl trans  drv     cty   hwy fl
##    <chr>        <chr>      <dbl> <int> <int> <chr>  <chr> <int> <int> <chr>
##  1 audi         a6 quatt~    4.2  2008     8 auto(~ 4        16    23 p
##  2 chevrolet    c1500 su~    5.3  2008     8 auto(~ r        14    20 r
##  3 chevrolet    c1500 su~    5.3  2008     8 auto(~ r        11    15 e
##  4 chevrolet    c1500 su~    5.3  2008     8 auto(~ r        14    20 r
##  5 chevrolet    c1500 su~    5.7  1999     8 auto(~ r        13    17 r
##  6 chevrolet    c1500 su~    6    2008     8 auto(~ r        12    17 r
##  7 chevrolet    corvette     5.7  1999     8 manua~ r        16    26 p
##  8 chevrolet    corvette     5.7  1999     8 auto(~ r        15    23 p
##  9 chevrolet    corvette     6.2  2008     8 manua~ r        16    26 p
## 10 chevrolet    corvette     6.2  2008     8 auto(~ r        15    25 p
## # ... with 60 more rows, and 1 more variable: class <chr>
```

**Section 5.2.4: #1, #3 and #4 only. You will need to install the nycflights13 package and use the flights data.**

# 1

Had an arrival delay of two or more hours

```
glimpse(flights)
```

```
## Observations: 336,776
## Variables: 19
## $ year          <int> 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013,...
## $ month         <int> 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,...
## $ dep_time      <int> 517, 533, 542, 544, 554, 554, 555, 557, 557, 55...
## $ sched_dep_time <int> 515, 529, 540, 545, 600, 558, 600, 600, 600, 60...
## $ dep_delay     <dbl> 2, 4, 2, -1, -6, -4, -5, -3, -3, -2, -2, -2, -2...
## $ arr_time      <int> 830, 850, 923, 1004, 812, 740, 913, 709, 838, 7...
## $ sched_arr_time <int> 819, 830, 850, 1022, 837, 728, 854, 723, 846, 7...
## $ arr_delay     <dbl> 11, 20, 33, -18, -25, 12, 19, -14, -8, 8, -2, -...
## $ carrier       <chr> "UA", "UA", "AA", "B6", "DL", "UA", "B6", "EV",...
## $ flight        <int> 1545, 1714, 1141, 725, 461, 1696, 507, 5708, 79...
## $ tailnum       <chr> "N14228", "N24211", "N619AA", "N804JB", "N668DN...
## $ origin        <chr> "EWR", "LGA", "JFK", "JFK", "LGA", "EWR", "EWR"...
## $ dest          <chr> "IAH", "IAH", "MIA", "BQN", "ATL", "ORD", "FLL"...
## $ air_time      <dbl> 227, 227, 160, 183, 116, 150, 158, 53, 140, 138...
## $ distance      <dbl> 1400, 1416, 1089, 1576, 762, 719, 1065, 229, 94...
## $ hour          <dbl> 5, 5, 5, 5, 6, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 5,...
## $ minute        <dbl> 15, 29, 40, 45, 0, 58, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ time_hour     <dttm> 2013-01-01 05:00:00, 2013-01-01 05:00:00, 2013...
```
```r
filter(flights, arr_delay > 120)
```

```
## # A tibble: 10,034 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      811            630       101     1047
## 2  2013     1     1      848           1835       853     1001
## 3  2013     1     1      957            733       144     1056
## 4  2013     1     1     1114            900       134     1447
## 5  2013     1     1     1505           1310       115     1638
## 6  2013     1     1     1525           1340       105     1831
## 7  2013     1     1     1549           1445        64     1912
## 8  2013     1     1     1558           1359       119     1718
## 9  2013     1     1     1732           1630        62     2028
## 10 2013     1     1     1803           1620       103     2008
## # ... with 10,024 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 <dttm>
```

Flew to Houston (IAH or HOU)

```r
filter(flights, dest %in% c("HOU","IAH"))
```

```
## # A tibble: 9,313 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      623            627        -4      933
## 4  2013     1     1      728            732        -4     1041
## 5  2013     1     1      739            739         0     1104
## 6  2013     1     1      908            908         0     1228
## 7  2013     1     1     1028           1026         2     1350
```

```
## 8  2013     1     1    1044          1045         -1    1352
## 9  2013     1     1    1114           900        134    1447
## 10 2013     1     1    1205          1200          5    1503
## # ... with 9,303 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 <dttm>
```

Were operated by United, American, or Delta

```
filter(flights, carrier %in% c("UA","AA","DL"))
```

```
## # A tibble: 139,504 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      554            600        -6      812
## 5  2013     1     1      554            558        -4      740
## 6  2013     1     1      558            600        -2      753
## 7  2013     1     1      558            600        -2      924
## 8  2013     1     1      558            600        -2      923
## 9  2013     1     1      559            600        -1      941
## 10 2013     1     1      559            600        -1      854
## # ... with 139,494 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 <dttm>
```

Departed in summer (July, August, and September)

```
filter(flights, month >= 7 ,  month <=9)
```

```
## # A tibble: 86,326 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>
## 1  2013     7     1        1           2029       212      236
## 2  2013     7     1        2           2359         3      344
## 3  2013     7     1       29           2245       104      151
## 4  2013     7     1       43           2130       193      322
## 5  2013     7     1       44           2150       174      300
## 6  2013     7     1       46           2051       235      304
## 7  2013     7     1       48           2001       287      308
## 8  2013     7     1       58           2155       183      335
## 9  2013     7     1      100           2146       194      327
## 10 2013     7     1      100           2245       135      337
## # ... with 86,316 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 <dttm>
```

Arrived more than two hours late, but didn't leave late

```
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
```

```
##      <int> <int> <int>    <int>          <int>    <dbl>    <int>
##  1   2013     1    27    1419           1420       -1     1754
##  2   2013    10     7    1350           1350        0     1736
##  3   2013    10     7    1357           1359       -2     1858
##  4   2013    10    16     657            700       -3     1258
##  5   2013    11     1     658            700       -2     1329
##  6   2013     3    18    1844           1847       -3       39
##  7   2013     4    17    1635           1640       -5     2049
##  8   2013     4    18     558            600       -2     1149
##  9   2013     4    18     655            700       -5     1213
## 10   2013     5    22    1827           1830       -3     2217
## # ... with 19 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 <dttm>
```

Were delayed by at least an hour, but made up over 30 minutes in flight

```
filter(flights,  dep_delay >=60, (dep_delay - arr_delay) > 30)
```

```
## # A tibble: 1,844 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    2205           1720      285       46
##  2   2013     1     1    2326           2130      116      131
##  3   2013     1     3    1503           1221      162     1803
##  4   2013     1     3    1839           1700       99     2056
##  5   2013     1     3    1850           1745       65     2148
##  6   2013     1     3    1941           1759      102     2246
##  7   2013     1     3    1950           1845       65     2228
##  8   2013     1     3    2015           1915       60     2135
##  9   2013     1     3    2257           2000      177       45
## 10   2013     1     4    1917           1700      137     2135
## # ... with 1,834 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 <dttm>
```

Departed between midnight and 6am (inclusive). Note that in dep_time, midnight is 2400, not 0.

```
filter(flights,  dep_time <=600 )
```

```
## # A tibble: 9,344 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 9,334 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 <dttm>
```

## 3

How many flights have a missing dep_time? What other variables are missing? What might these rows represent?

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

```
## # A tibble: 8,255 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       NA           1630        NA       NA
## 2   2013     1     1       NA           1935        NA       NA
## 3   2013     1     1       NA           1500        NA       NA
## 4   2013     1     1       NA            600        NA       NA
## 5   2013     1     2       NA           1540        NA       NA
## 6   2013     1     2       NA           1620        NA       NA
## 7   2013     1     2       NA           1355        NA       NA
## 8   2013     1     2       NA           1420        NA       NA
## 9   2013     1     2       NA           1321        NA       NA
## 10  2013     1     2       NA           1545        NA       NA
## # ... with 8,245 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 <dttm>
```

These flights never took off, so cancelled. #4 Why is NA ^ 0 not missing? Why is NA | TRUE not missing? Why is FALSE & NA not missing? Can you figure out the general rule? (NA * 0 is a tricky counterexample!)

**Need to check with Robert.**

**Section 5.4.1: #1 and #3 only**

# 1 Brainstorm as many ways as possible to select dep_time, dep_delay, arr_time, and arr_delay from flights.

```
select(flights,dep_time,dep_delay,arr_time,arr_delay)
```

```
## # A tibble: 336,776 x 4
##    dep_time dep_delay arr_time arr_delay
##       <int>     <dbl>    <int>     <dbl>
## 1       517         2      830        11
## 2       533         4      850        20
## 3       542         2      923        33
## 4       544        -1     1004       -18
## 5       554        -6      812       -25
## 6       554        -4      740        12
## 7       555        -5      913        19
```

```
## 8       557       -3       709      -14
## 9       557       -3       838       -8
## 10      558       -2       753        8
## # ... with 336,766 more rows
```

## 2 What happens if you include the name of a variable multiple times in a select() call?

```
select(flights,dep_time,dep_time,arr_time,arr_delay)
```

```
## # A tibble: 336,776 x 3
##    dep_time arr_time arr_delay
##       <int>    <int>     <dbl>
## 1       517      830        11
## 2       533      850        20
## 3       542      923        33
## 4       544     1004       -18
## 5       554      812       -25
## 6       554      740        12
## 7       555      913        19
## 8       557      709       -14
## 9       557      838        -8
## 10      558      753         8
## # ... with 336,766 more rows
```

R Markdown will remove the duplicate variable name and gives only variables.

## 3 What does the one_of() function do? Why might it be helpful in conjunction with this vector?

## 4 Does the result of running the following code surprise you? How do the select helpers deal with case by default? How can you change that default?

```
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 <dttm>
```