

# MA615\_Assignment2

*Siwei Hu*

*September 22, 2018*

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.4.4
```

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

```
## v ggplot2 3.0.0      v purrr  0.2.5
```

```
## v tibble  1.4.2      v dplyr  0.7.6
```

```
## v tidyr   0.8.1      v stringr 1.2.0
```

```
## v readr   1.1.1      v forcats 0.3.0
```

```
## Warning: package 'ggplot2' was built under R version 3.4.4
```

```
## Warning: package 'tibble' was built under R version 3.4.4
```

```
## Warning: package 'tidyr' was built under R version 3.4.4
```

```
## Warning: package 'readr' was built under R version 3.4.4
```

```
## Warning: package 'purrr' was built under R version 3.4.4
```

```
## Warning: package 'dplyr' was built under R version 3.4.4
```

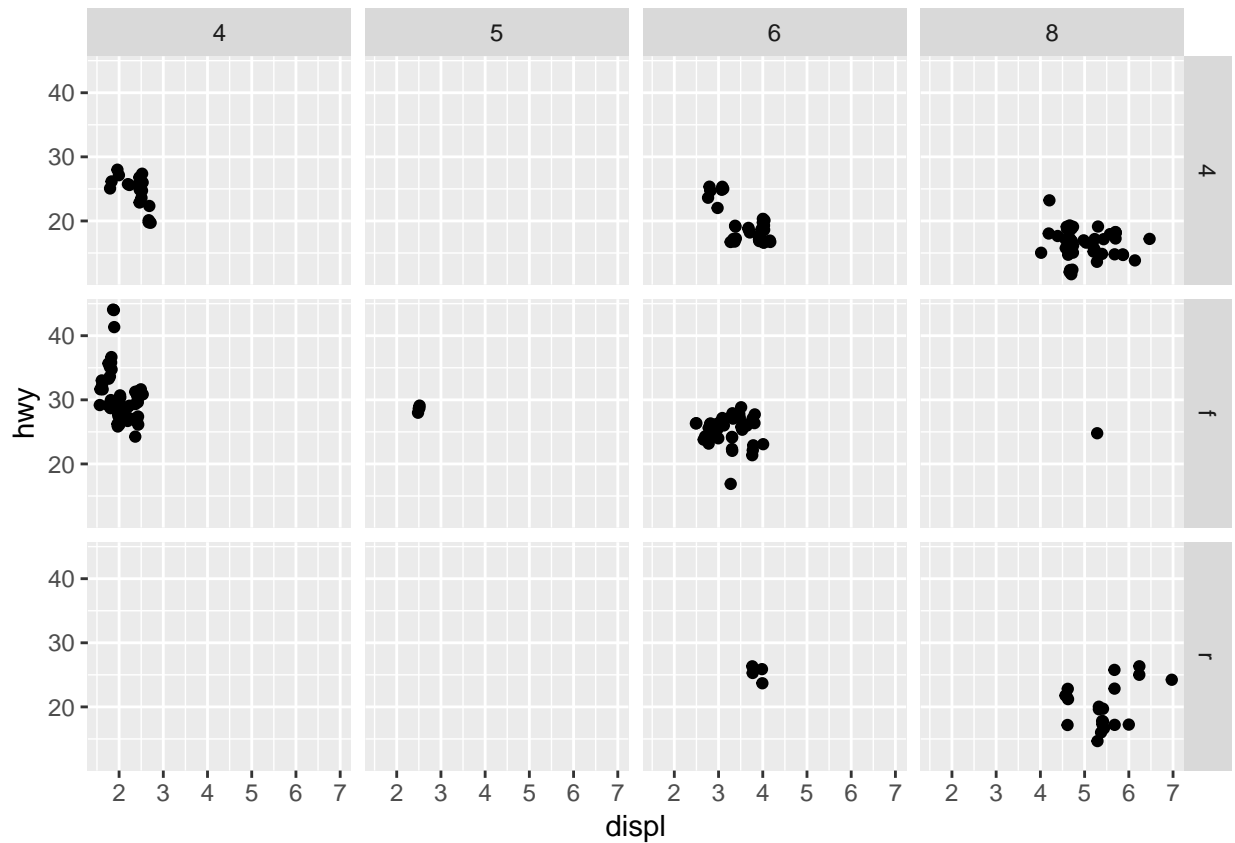
```
## Warning: package 'forcats' was built under R version 3.4.4
```

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

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

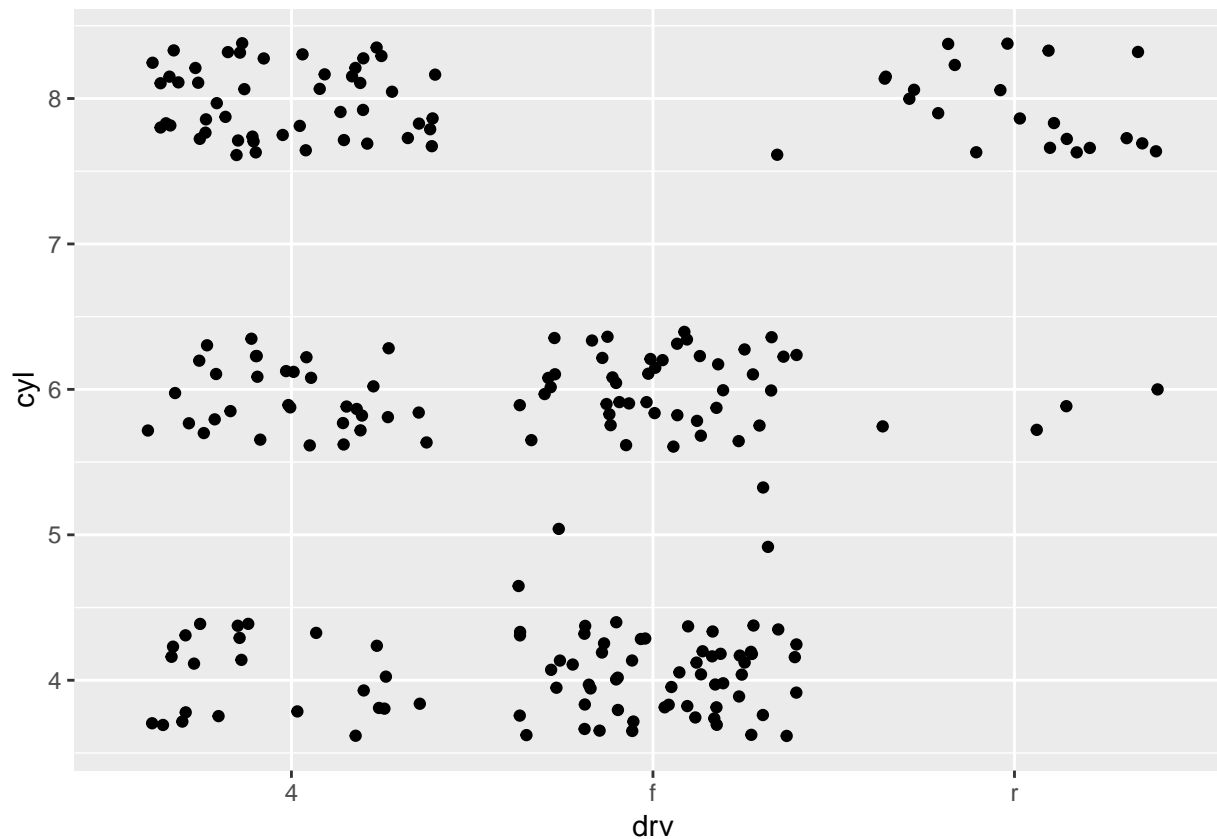
```
## x dplyr::lag()     masks stats::lag()
```

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



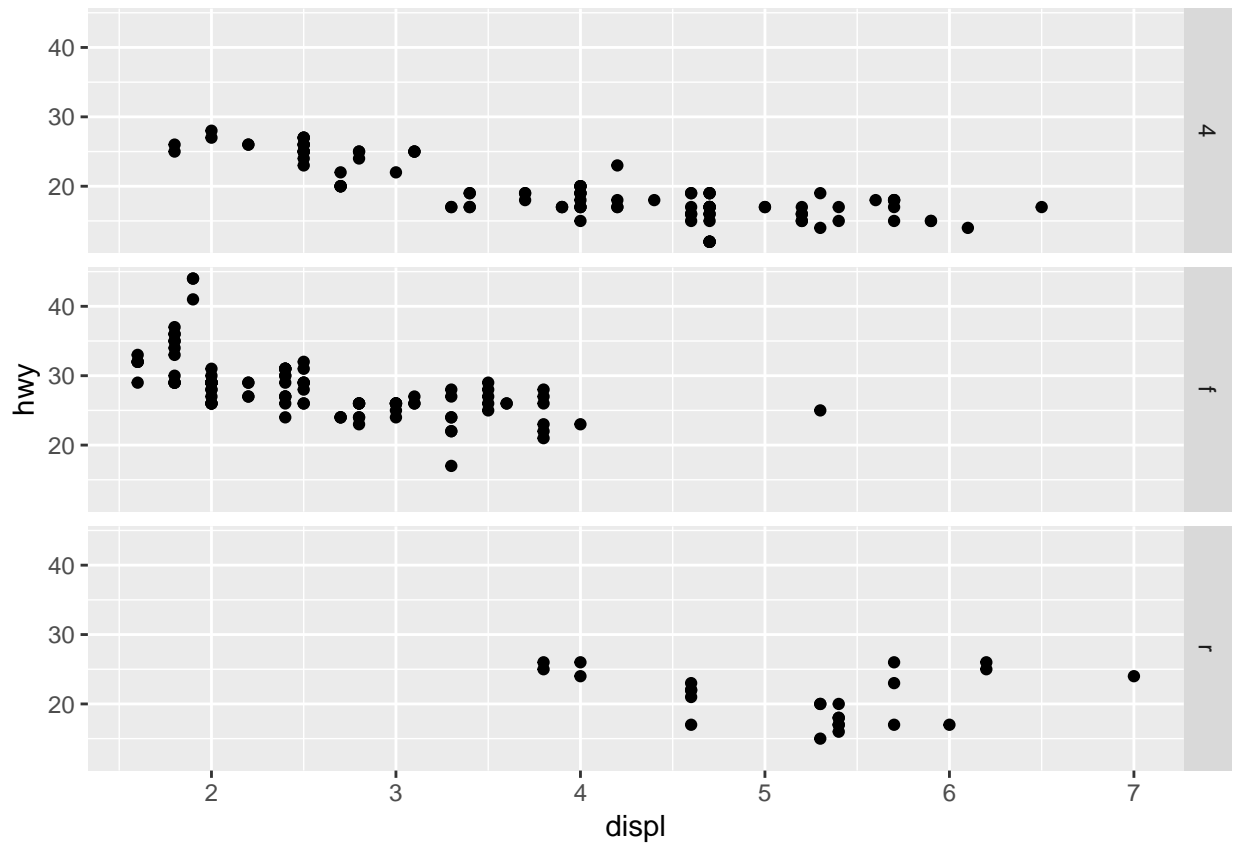
The empty cell help the plot of linear regression between displ and hwy to be categorized with two different variables drv and cyl. Both variables are binary. So the plot changes to 12 different pieces with 12 different situations. In each piece, they will do linear regression between displ and hwy again.

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

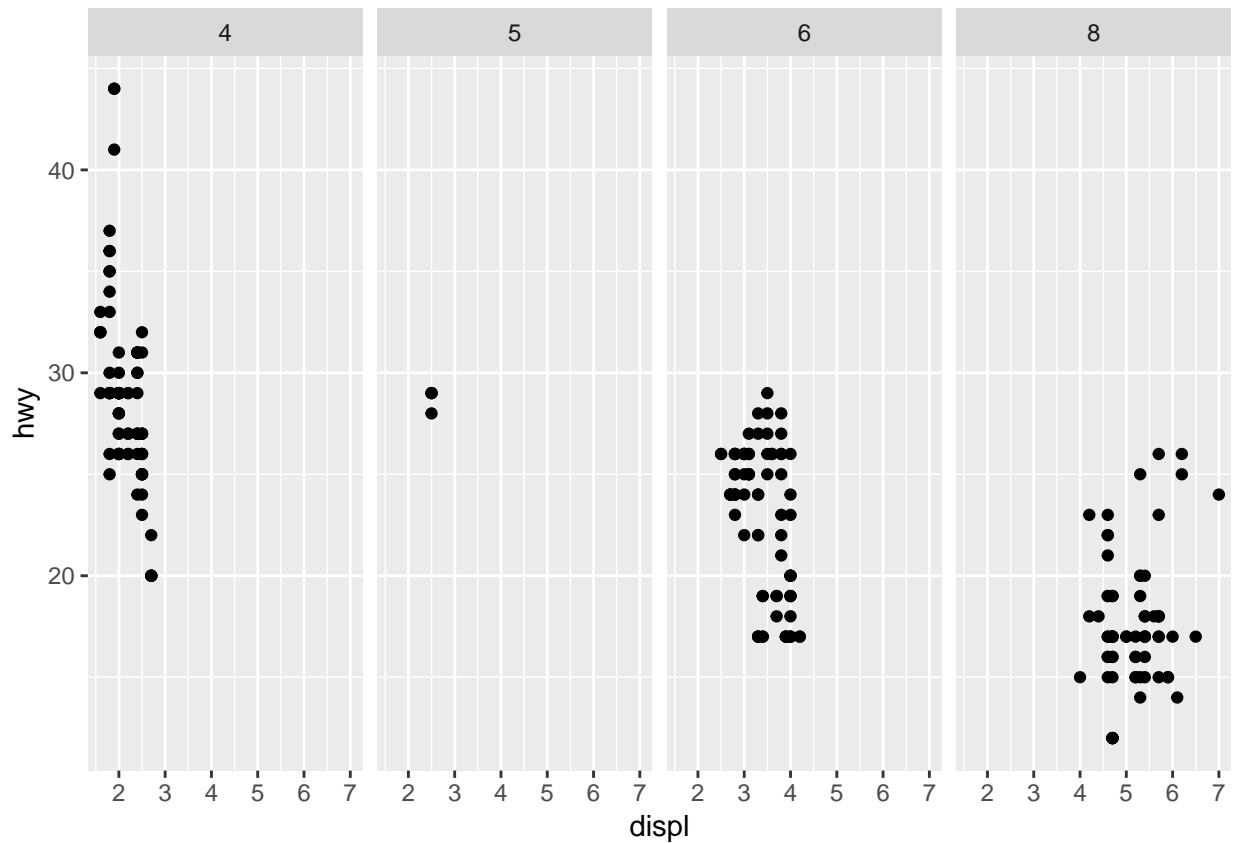


This plot categorized all car and we can see they are in 12 different situation. This is same as the previous graph. The differences between them are that previous graph did linear regression between displ and hwy in each piece. The number of points in each piece is same.

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



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

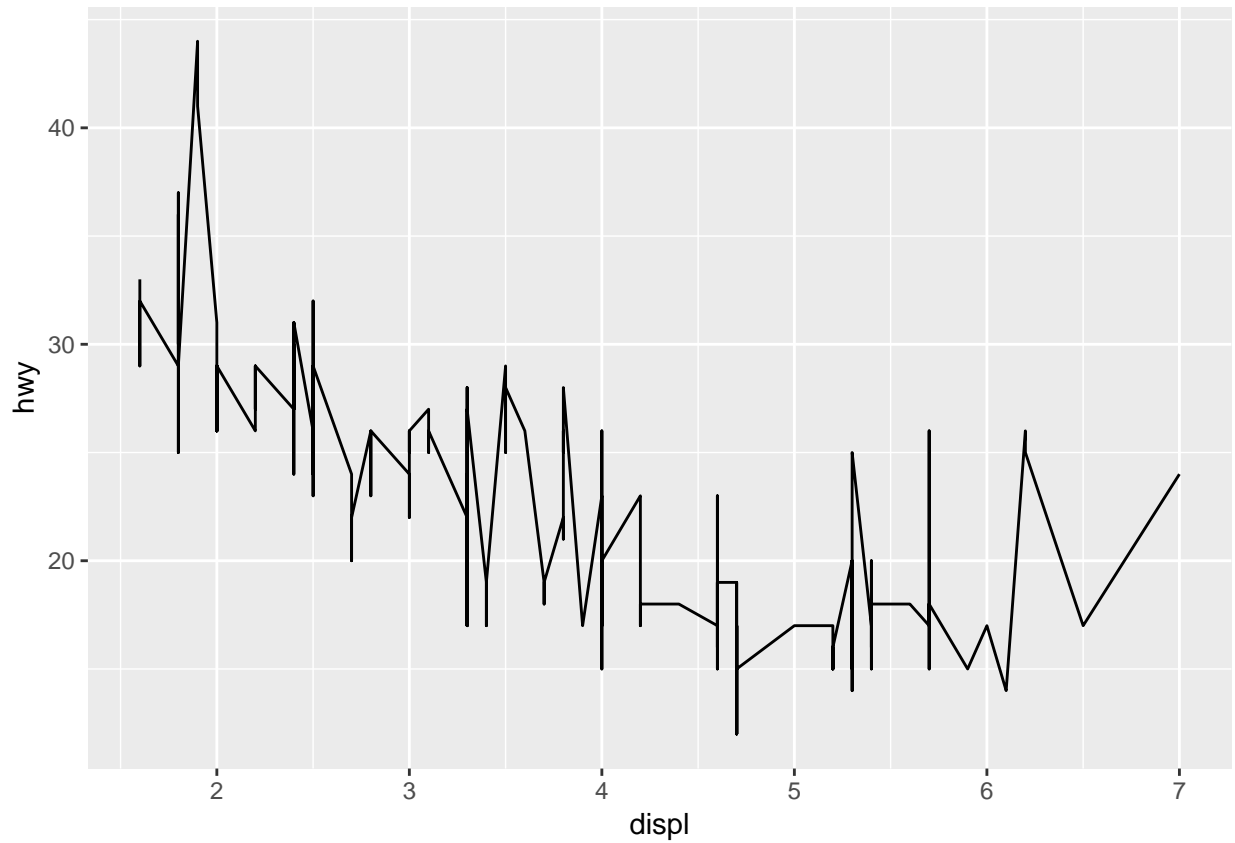


The first graph divide data into three rows of graph that use variable “drv”. So there are three rows and it does linear regression between displ and hwy in each row.

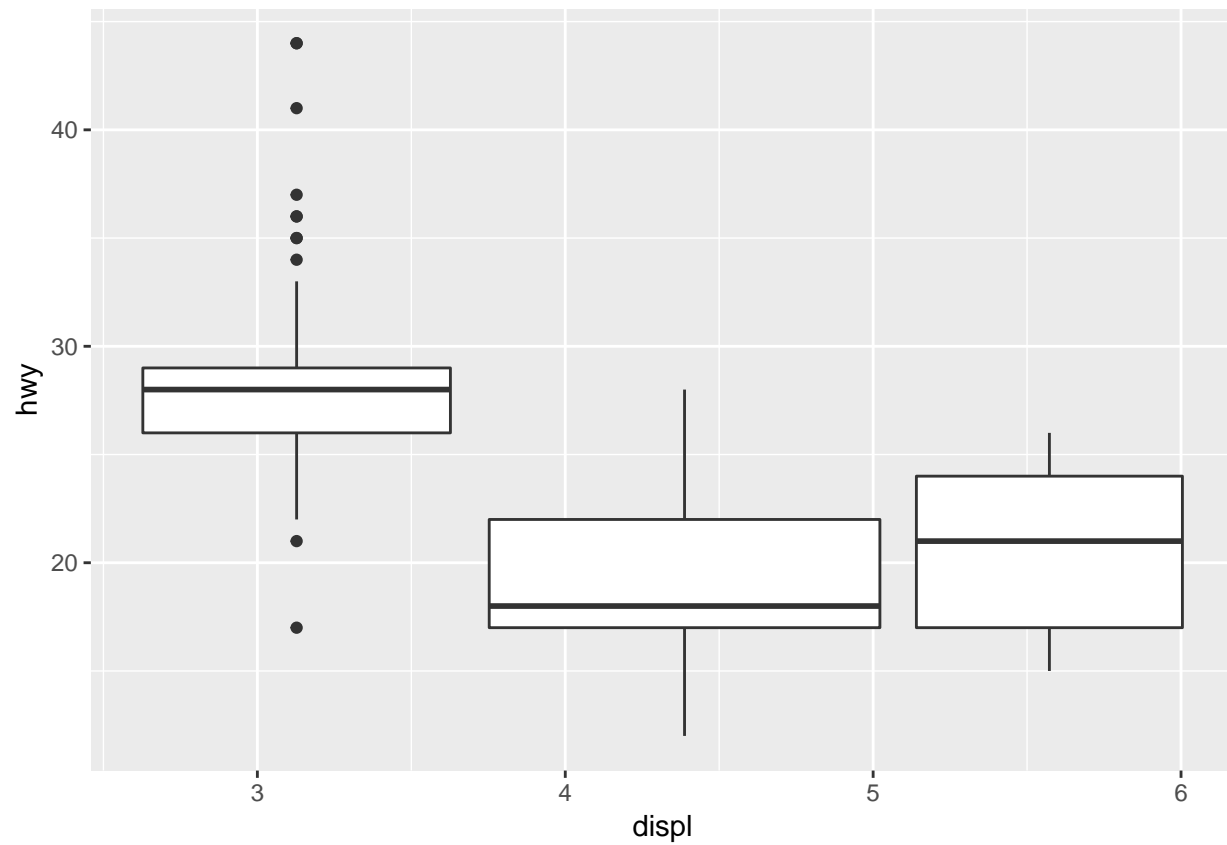
The Second graph divide data into four columns of graph that use variable “cyl”. So there are four columns and it does linear regression between displ and hwy in each column.

The . means if you dont want to do facet in row or column, it can instead ofthat variable name.

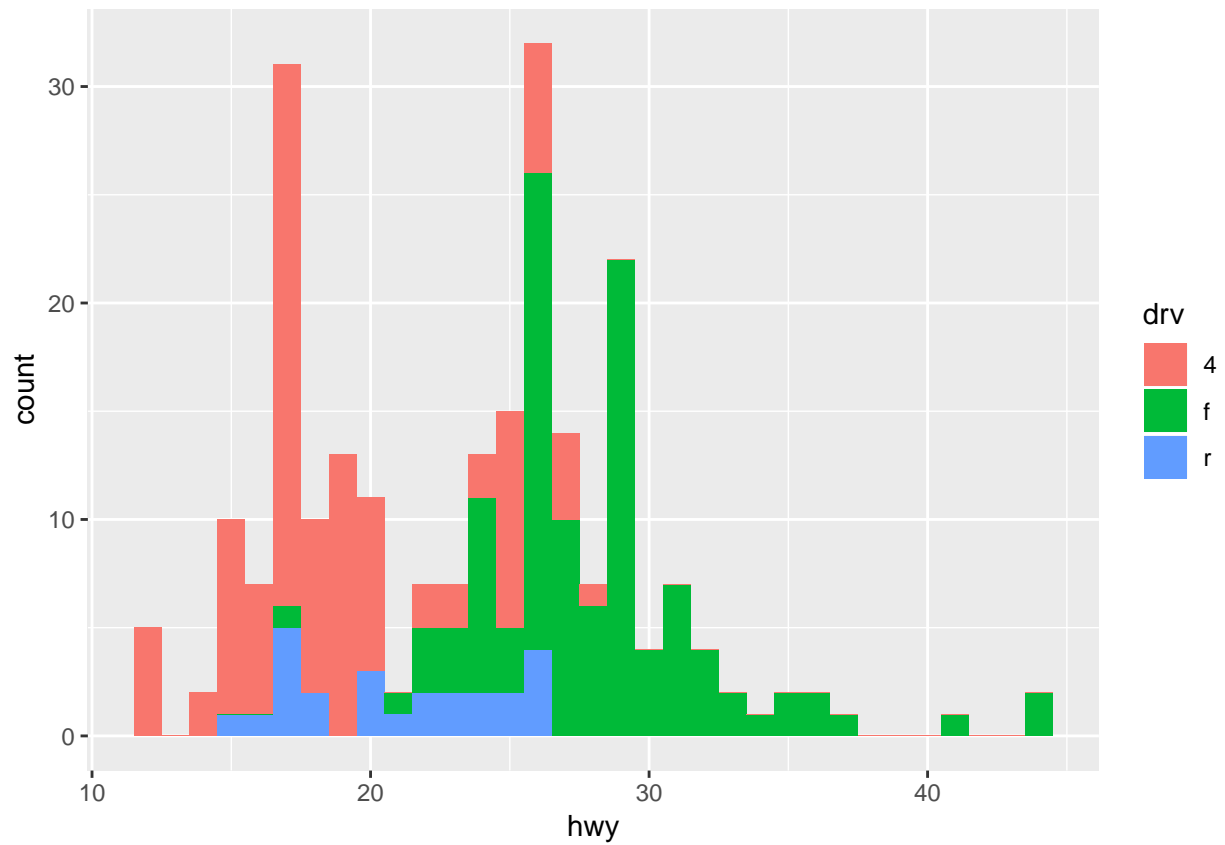
```
ggplot(data = mpg)+
  geom_line(mapping = aes(x = displ,y = hwy))
```



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

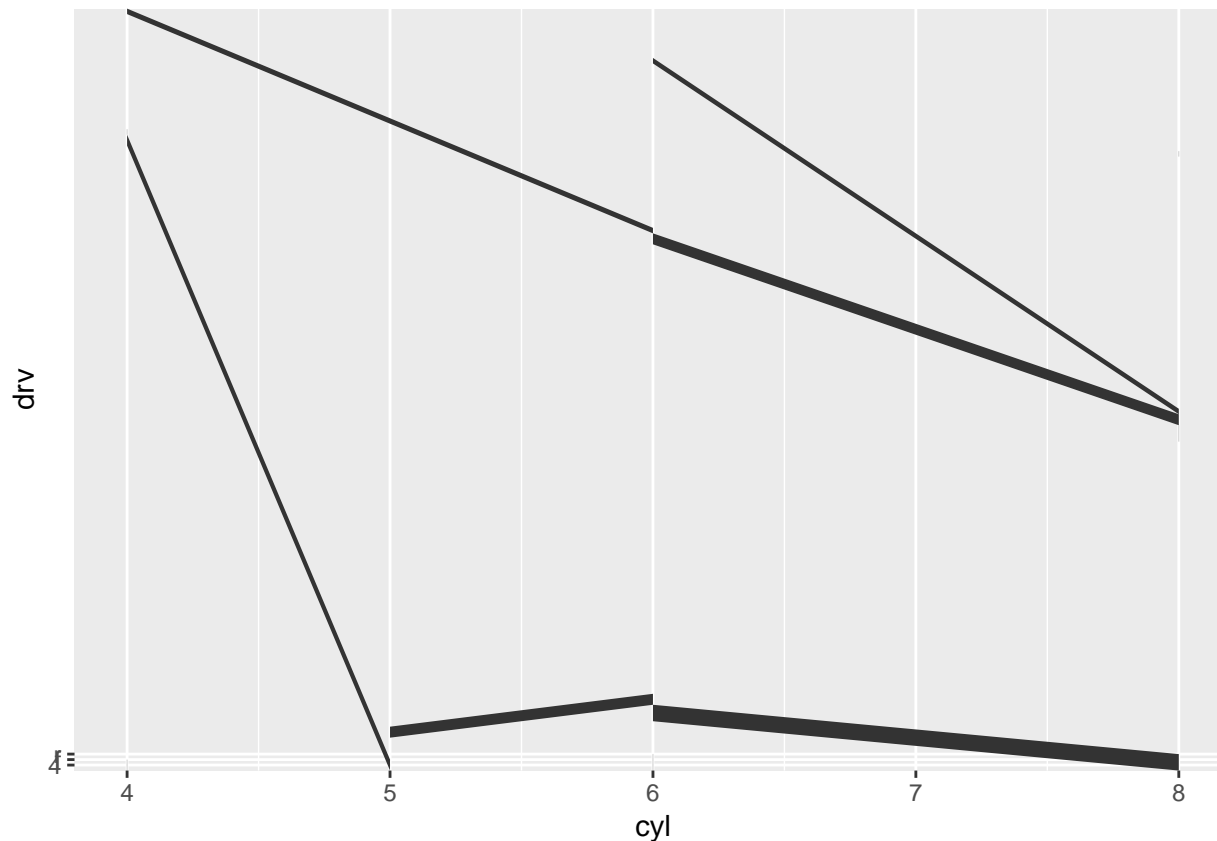


```
ggplot(data = mpg)+  
  geom_histogram(mapping = aes(x = hwy, fill = drv), binwidth = 1)
```



```
ggplot(data = mpg)+  
  geom_area(mapping = aes(x = cyl,y = drv))
```





```
library(nycflights13)
```

```
## Warning: package 'nycflights13' was built under R version 3.4.4
```

```
nycflights13::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>
```

```
filter(flights, arr_delay >= 2)
```

```
## Warning: package 'bindrcpp' was built under R version 3.4.4
```

```
## # A tibble: 127,929 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             558          -4     740
## 5  2013     1     1     555             600          -5     913
## 6  2013     1     1     558             600          -2     753
## 7  2013     1     1     558             600          -2     924
## 8  2013     1     1     559             600          -1     941
## 9  2013     1     1     600             600           0     837
##10  2013     1     1     602             605          -3     821
## # ... with 127,919 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, dest %in% c('IAH', 'HOU'))
```

```
## # 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 <dtm>
```

```
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 <dtm>
```

```
filter(flights, month %in% c(7,8,9))
```

```
## # 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 <dtm>
```

```
filter(flights, arr_delay >= 2 & dep_delay == 0 )
```

```
## # A tibble: 4,694 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     600           600           0     837
## 2  2013     1     1     635           635           0    1028
## 3  2013     1     1     739           739           0    1104
## 4  2013     1     1     745           745           0    1135
## 5  2013     1     1     800           800           0    1022
## 6  2013     1     1     805           805           0    1015
## 7  2013     1     1     810           810           0    1048
## 8  2013     1     1     823           823           0    1151
## 9  2013     1     1     830           830           0    1018
## 10 2013     1     1     835           835           0    1210
```

```
## # ... with 4,684 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, dep_delay >= 1 & dep_time - sched_dep_time >= 30 )
```

```
## # A tibble: 62,097 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     732           645          47    1011
## 2  2013     1     1     749           710          39     939
## 3  2013     1     1     811           630         101    1047
## 4  2013     1     1     826           715          71    1136
## 5  2013     1     1     903           820          43    1045
## 6  2013     1     1     906           843          23    1134
## 7  2013     1     1     909           810          59    1331
## 8  2013     1     1     953           921          32    1320
```

```
## 9 2013 1 1 957 733 144 1056
## 10 2013 1 1 1003 959 4 1408
## # ... with 62,087 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, dep_time >= 0 & 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 <dtm>
```

```
filter(flights, between(dep_time,0,600))
```

```
## # 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 <dtm>
```

Between(x, left, right) This help me easier to set a range of time in preivous quesiton.

```
miss.dep <- sum(length(which(is.na(flights$dep_time))))
miss.dep
```

```
## [1] 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
##   <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 <dtm>
```

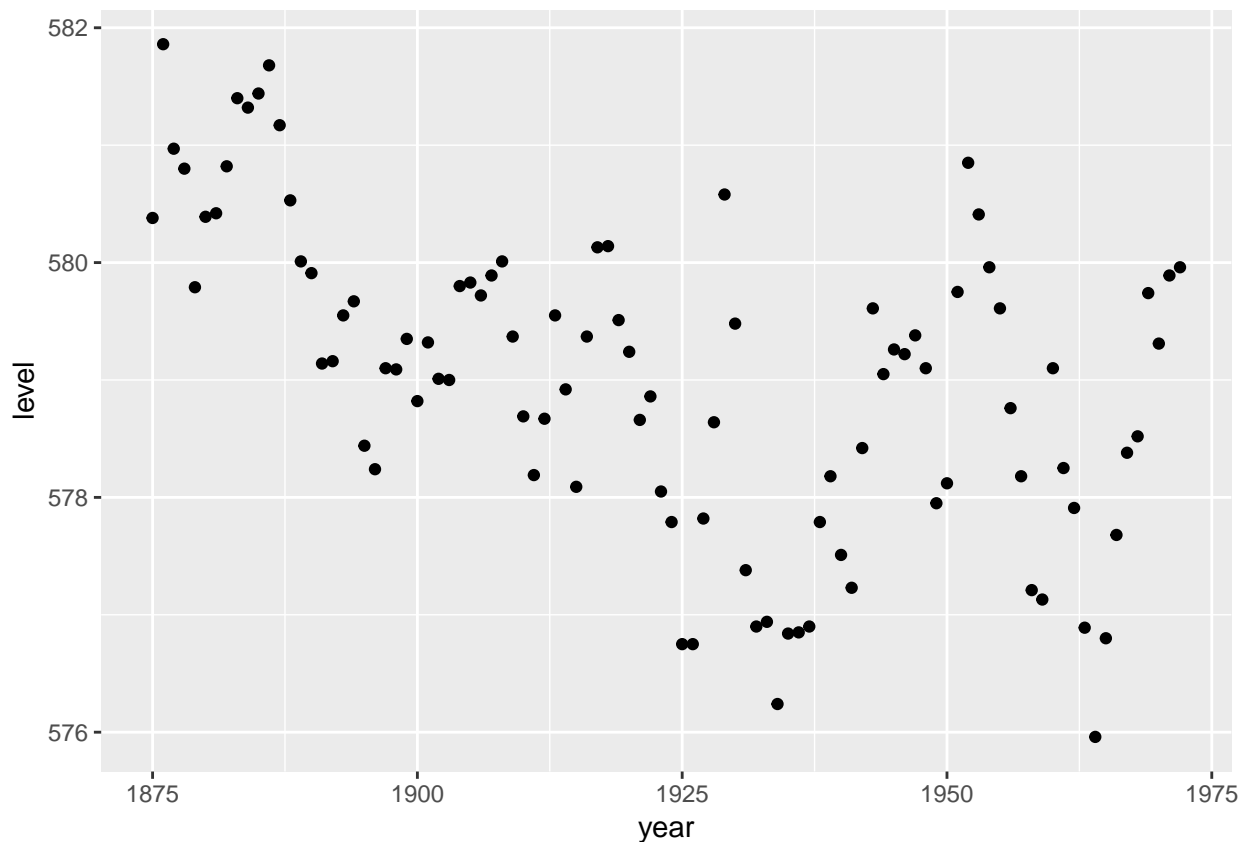
If there is no departure time, there are no arr\_time. It means the flights were canceled.

NA^0, NA|TRUE and FALSE&NA are be

```
data("LakeHuron")
LakeHuron <- data.frame("year" = 1875:1972, "level" = LakeHuron)

ggplot(data = LakeHuron)+
  geom_point(mapping = aes(x = year,y = level))
```

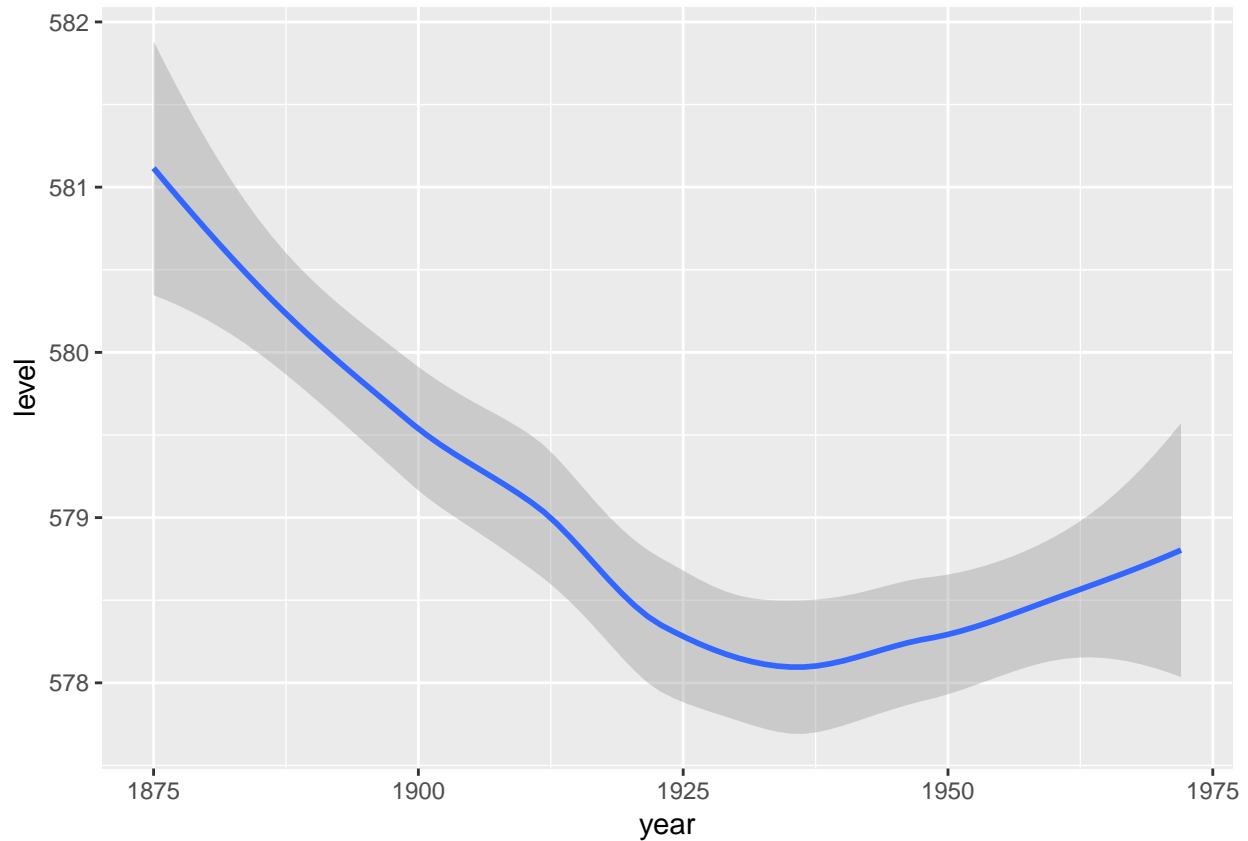
## Don't know how to automatically pick scale for object of type ts. Defaulting to continuous.



```
ggplot(data = LakeHuron)+
  geom_smooth(mapping = aes(x = year,y = level))
```

## Don't know how to automatically pick scale for object of type ts. Defaulting to continuous.

## `geom\_smooth()` using method = 'loess' and formula 'y ~ x'



```
ggplot(data = LakeHuron)+
  geom_point(mapping = aes(x = year,y = level))+
  geom_smooth(mapping = aes(x= year, y= level),se = FALSE)
```

## Don't know how to automatically pick scale for object of type ts. Defaulting to continuous.

## `geom\_smooth()` using method = 'loess' and formula 'y ~ x'

