

Analyze Boston

Michael Rose

```
knitr::opts_chunk$set(echo = TRUE)
```

Abstract

‘Analyze Boston’ is an open data initiative maintained by the city of Boston containing facts, figures, and maps related to the city. In this project we will look at some of the city’s 133+ data sets and analyze them with descriptive and inferential statistics. The focus of the project is the discovery of interesting patterns.

Intro to the Data

In this analysis, I will be using 4 different data sets.

Employee Earnings Report

- Each year, the City of Boston publishes payroll data for employees. This dataset contains employee names, job details, and earnings information including base salary, overtime, and total compensation for employees of the City.
- You can see more at <https://data.boston.gov/dataset/employee-earnings-report>

```
earnings_2017 %>% head(n = 10L)
```

```
## # A tibble: 22,235 x 12
##   NAME      `DEPARTMENT NAME` TITLE  REGULAR RETRO  OTHER OVERTIME INJURED
##   <chr>      <chr>              <chr>  <dbl> <dbl>  <dbl>  <dbl>  <dbl>
## 1 Miller~ Boston Police De~ Police~ 129531.  NA  13694.  8150.  NA
## 2 Sulliv~ Boston Police De~ Office~ 56922.  NA  3595.  1548.  NA
## 3 O'Hara~ Boston Police De~ Police~ 124057.  NA  6432.  29044.  NA
## 4 Whalen~ Boston Police De~ Police~ 94956.  4985. 13592.  85419.  58.0
## 5 Kelly,~ Boston Police De~ Tape L~ 69995.  NA  300  7961.  NA
## 6 Carrol~ Boston Police De~ Police~ 12757.  2390. 41612.  912.  NA
## 7 Connol~ Boston Police De~ Police~ 93180.  2028. 13338.  19882.  NA
## 8 Ivens,~ Boston Police De~ Police~  NA  NA  60777.  NA  2659.
## 9 Kelly,~ Boston Police De~ Police~ 13827.  NA  62393.  868.  NA
## 10 Klokma~ Boston Police De~ Police~ 107599.  NA  14482.  12825.  NA
## # ... with 22,225 more rows, and 4 more variables: DETAIL <dbl>,
## # `QUINN/EDUCATION INCENTIVE` <dbl>, `TOTAL EARNINGS` <dbl>,
## # POSTAL <int>
```

Crime Incident Reports

- Crime incident reports are provided by Boston Police Department (BPD) to document the initial details surrounding an incident to which BPD officers respond. This is a dataset containing records from the new crime incident report system, which includes a reduced set of fields focused on capturing the type of incident as well as when and where it occurred. Records in the new system begin in June of 2015.

- You can see more at <https://data.boston.gov/dataset/crime-incident-reports-august-2015-to-date-source-new-system>

```
crime %>% head()
```

```
## # A tibble: 6 x 17
##   INCIDENT_NUMBER OFFENSE_CODE OFFENSE_CODE_GRO~ OFFENSE_DESCRIP~ DISTRICT
##   <chr>           <chr>           <chr>           <chr>           <chr>
## 1 I182017604      03115           Investigate Pers~ INVESTIGATE PER~ B3
## 2 I182017601      00520           Residential Burg~ BURGLARY - RESI~ B2
## 3 I182017596      03831           Motor Vehicle Ac~ M/V - LEAVING S~ A1
## 4 I182017595      03802           Motor Vehicle Ac~ "M/V ACCIDENT -- <NA>
## 5 I182017594      01830           Drug Violation    DRUGS - SICK AS~ D14
## 6 I182017593      00361           Robbery           ROBBERY - OTHER  D4
## # ... with 12 more variables: REPORTING_AREA <int>, SHOOTING <chr>,
## #   OCCURRED_ON_DATE <dtm>, YEAR <int>, MONTH <int>, DAY_OF_WEEK <chr>,
## #   HOUR <int>, UCR_PART <chr>, STREET <chr>, Lat <dbl>, Long <dbl>,
## #   Location <chr>
```

BPD Firearm Recovery Counts

- This dataset provides daily counts of firearms recovered by Boston Police Department since August 20, 2014. Recovery totals are provided for three distinct channels: crime, voluntary surrender, and gun buyback programs.
- You can see more at <https://data.boston.gov/dataset/boston-police-department-firearms-recovery-counts>

```
gun_recovery %>% head()
```

```
## # A tibble: 6 x 4
##   CollectionDate CrimeGunsRecovered GunsSurrenderedSafe~ BuybackGunsRecov~
##   <chr>           <int>           <int>           <int>
## 1 8/20/2014      2              3              1
## 2 8/21/2014      2              0              4
## 3 8/22/2014      0              0              2
## 4 8/25/2014      8              3              0
## 5 8/26/2014      9              0              0
## 6 8/27/2014      1              0              0
```

Economic Indicators

- The Boston Planning and Redevelopment Authority (BPDA), formerly known as the Boston Redevelopment Authority (BRA), is tasked with planning for and guiding inclusive growth within the City of Boston. As part of this mission, BPDA collects and analyzes a variety of economic data relating to topics such as the employment, housing, travel, and real estate development. This is a legacy dataset of economic indicators tracked monthly between January 2013 and January 2015.
- You can see more at <https://data.boston.gov/dataset/economic-indicators-legacy-portal>

```
econ_indicators %>% head()
```

```
## # A tibble: 6 x 19
##   Year Month logan_passengers logan_intl_flights hotel_occup_rate
##   <int> <int>           <int>           <int>           <dbl>
## 1 2013     1      2019662          2986          0.572
## 2 2013     2      1878731          2587          0.645
## 3 2013     3      2469155          3250          0.819
```

```

## 4 2013      4      2551246      3408      0.855
## 5 2013      5      2676291      3240      0.858
## 6 2013      6      2824862      3402      0.911
## # ... with 14 more variables: hotel_avg_daily_rate <dbl>,
## #   total_jobs <int>, unemp_rate <dbl>, labor_force_part_rate <dbl>,
## #   pipeline_unit <int>, pipeline_total_dev_cost <dbl>,
## #   pipeline_sqft <int>, pipeline_const_jobs <dbl>, foreclosure_pet <int>,
## #   foreclosure_deeds <int>, med_housing_price <int>,
## #   housing_sales_vol <int>, new_housing_const_permits <int>,
## #   `new-affordable_housing_permits` <int>

```

Data Cleaning

The first thing that needs to be done is tidying up the data. We can start by removing any numeric NAs and turning them into 0s.

```
# check to see if there are any NA entries in non numeric  
# columns
```

```
sum(is.na(earnings_2017$NAME))
```

```
## [1] 0
```

```
sum(is.na(earnings_2017$"DEPARTMENT NAME"))
```

```
## [1] 0
```

```
sum(is.na(earnings_2017$TITLE))
```

```
## [1] 0
```

```
# we are good!
```

```
# set NA to 0 in numeric columns
```

```
earnings_2017[is.na(earnings_2017)] <- 0
```

Since minimum wage is \$11/hr we can filter full time from part time. I will be removing those who make under $11 * 40 * 52 = \$22880$ / year. I will round down to \$20,000.

```
# We can also do some mutating to give us a cleaner data  
# frame
```

```
earnings_2017 <- earnings_2017 %>% filter(earnings_2017$`TOTAL EARNINGS` >  
  20000) %>% mutate(EXTRA_PAY = .$RETRO + .$OTHER + .$DETAIL +  
  .$`QUINN/EDUCATION INCENTIVE`) %>% select("NAME", "DEPARTMENT NAME",  
  "TITLE", "REGULAR", "OVERTIME", "EXTRA_PAY", "TOTAL EARNINGS",  
  "INJURED")
```

```
earnings_2017 %>% head()
```

```
## # A tibble: 6 x 8
```

```
##   NAME `DEPARTMENT NAM~ TITLE REGULAR OVERTIME EXTRA_PAY `TOTAL EARNINGS`  
##   <chr> <chr>           <chr>   <dbl>   <dbl>   <dbl>         <dbl>  
## 1 Mill~ Boston Police D~ Poli~ 129531.   8150.   37981.   175663.  
## 2 Sull~ Boston Police D~ Offi~  56922.   1548.    3595.    62065.  
## 3 O'Ha~ Boston Police D~ Poli~ 124057.  29044.   52078.   205178.  
## 4 Whal~ Boston Police D~ Poli~  94956.  85419.   54878.   235312.  
## 5 Kell~ Boston Police D~ Tape~  69995.   7961.     300     78256.  
## 6 Carr~ Boston Police D~ Poli~  12757.    912.   45566.   59234.
```

```
## # ... with 1 more variable: INJURED <dbl>
```

Now we can begin to clean up those department name factors to get a clearer view of the groups as a whole. It seems like the school system takes up the majority of the factors, so lets compress them all into one factor - "Education"

```
# We want to combine all of the factors that are within the  
# education system into one education factor
```

```
# We start with 227 factors
```

```

head(factor(unique(earnings_2017$`DEPARTMENT NAME`)))

## [1] Boston Police Department      Workers Compensation Service
## [3] BPS East Boston High           BPS School Safety Service
## [5] Dpt of Innovation & Technology BPS Ohrenberger Elementary
## 224 Levels: Accountability Achievement Gap ... Youth Engagement & Employment
# find all factors containing BPS

earnings_2017_bps_factors <- str_detect(earnings_2017$`DEPARTMENT NAME`,
  "BPS")

head(unique(earnings_2017$`DEPARTMENT NAME`[earnings_2017_bps_factors]))

## [1] "BPS East Boston High"      "BPS School Safety Service"
## [3] "BPS Ohrenberger Elementary" "BPS Transportation"
## [5] "BPS Quincy Elementary"     "BPS McCormack Middle"
# now lets add K-8

earnings_2017_k8_factors <- str_detect(earnings_2017$`DEPARTMENT NAME`,
  "K-8")

head(unique(earnings_2017$`DEPARTMENT NAME`[earnings_2017_k8_factors]))

## [1] "Roosevelt K-8"      "Edison K-8"      "Higginson/Lewis K-8"
## [4] "Jackson/Mann K-8"   "Greenwood, S K-8" "Lyon K-8"
# Academy

earnings_2017_Academy_factors <- str_detect(earnings_2017$`DEPARTMENT NAME`,
  "Academy")

head(unique(earnings_2017$`DEPARTMENT NAME`[earnings_2017_Academy_factors]))

## [1] "Tech Boston Academy"      "BPS Latin Academy"
## [3] "West Roxbury Academy"     "BPS MPH\\Crafts Academy"
## [5] "Kennedy, EM Health Academy" "WREC: Urban Science Academy"
# Others

earnings_2017_otherEdu_factors <- c("Mattapan Early Elementary",
  "English Language Learn", "Kennedy, JF Elementary", "East Boston EEC",
  "Baldwin ELC", "Superintendent", "Unified Student Svc", "Umana Middle",
  "Achievement Gap", "Early Learning Services", "Kennedy, PJ Elementary",
  "Career & Technical Ed", "HPEC: Com Acd Science & Health",
  "Henderson Elementary", "Frederick Pilot Middle", "West Zone ELC",
  "Boston Collaborative High Sch", "Haley Pilot", "Enrollment Services",
  "Haynes EEC", "Teaching & Learning", "BTU Pilot", "Asst Superintendent-Network A",
  "Ellison/Parks EES", "Alighieri Montessori School", "Lyon Pilot High 9-12",
  "Dudley St. Neighborhood School", "P. A. Shaw Elementary",
  "UP \\Unlocking Potential\\ Acad", "Greater Egleston High",
  "Diplomas Plus", "Quincy Upper School", "Student Support Svc",
  "Chief Academic Officer")

# mutate

```

```

earnings_2017$`DEPARTMENT NAME`[earnings_2017_bps_factors] <- "Education"

head(unique(earnings_2017$`DEPARTMENT NAME`))

## [1] "Boston Police Department"      "Workers Compensation Service"
## [3] "Education"                    "Dpt of Innovation & Technology"
## [5] "Registry Division"            "Boston Fire Department"

# Down to 127 factors

earnings_2017$`DEPARTMENT NAME`[earnings_2017_k8_factors] <- "Education"

head(unique(earnings_2017$`DEPARTMENT NAME`))

## [1] "Boston Police Department"      "Workers Compensation Service"
## [3] "Education"                    "Dpt of Innovation & Technology"
## [5] "Registry Division"            "Boston Fire Department"

# Down to 108 factors

earnings_2017$`DEPARTMENT NAME`[earnings_2017_Academy_factors] <- "Education"

head(unique(earnings_2017$`DEPARTMENT NAME`))

## [1] "Boston Police Department"      "Workers Compensation Service"
## [3] "Education"                    "Dpt of Innovation & Technology"
## [5] "Registry Division"            "Boston Fire Department"

# Down to 97 factors

for (i in earnings_2017_otherEdu_factors) {
  logic <- str_detect(earnings_2017$`DEPARTMENT NAME`, i)
  earnings_2017$`DEPARTMENT NAME`[logic] <- "Education"
}

# Down to 63 factors

head(unique(earnings_2017$`DEPARTMENT NAME`))

## [1] "Boston Police Department"      "Workers Compensation Service"
## [3] "Education"                    "Dpt of Innovation & Technology"
## [5] "Registry Division"            "Boston Fire Department"

```

Education Pay

```
earnings_2017 %>% filter(`DEPARTMENT NAME` == "Education") %>%  
  arrange(desc(`TOTAL EARNINGS`)) %>% head()
```

```
## # A tibble: 6 x 8  
##   NAME `DEPARTMENT NAME` TITLE REGULAR OVERTIME EXTRA_PAY `TOTAL EARNINGS`  
##   <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl>  
## 1 Bott~ Education BPS ~ 0 0 0 285459.  
## 2 Chan~ Education Supe~ 264661. 0 6000. 270661.  
## 3 McCa~ Education Teac~ 46981. 0 182189. 229170.  
## 4 Jord~ Education Unit~ 106762. 0 81789. 188551.  
## 5 Estr~ Education Depu~ 177625. 0 0 177625.  
## 6 Wood~ Education Asst~ 4620 0 169047. 173667.  
## # ... with 1 more variable: INJURED <dbl>
```

From the above the pays seem about normal. The first woman Torii Bottomley won a lawsuit against her employer for workplace bullying:

https://www.pacermonitor.com/public/case/22846991/Bottomley_v_Boston_Public_Schools_et_al

The only other outlier that it shown is Elaine M McCabe with a base salary of \$46,981 and \$182,189.26 in extra pay (not overtime). I was unable to find a reasoning for this.

Police Pay

First lets look at the how the pay is distributed across the police department.

```
# Find average pay, sd of pay for police officers and plot
```

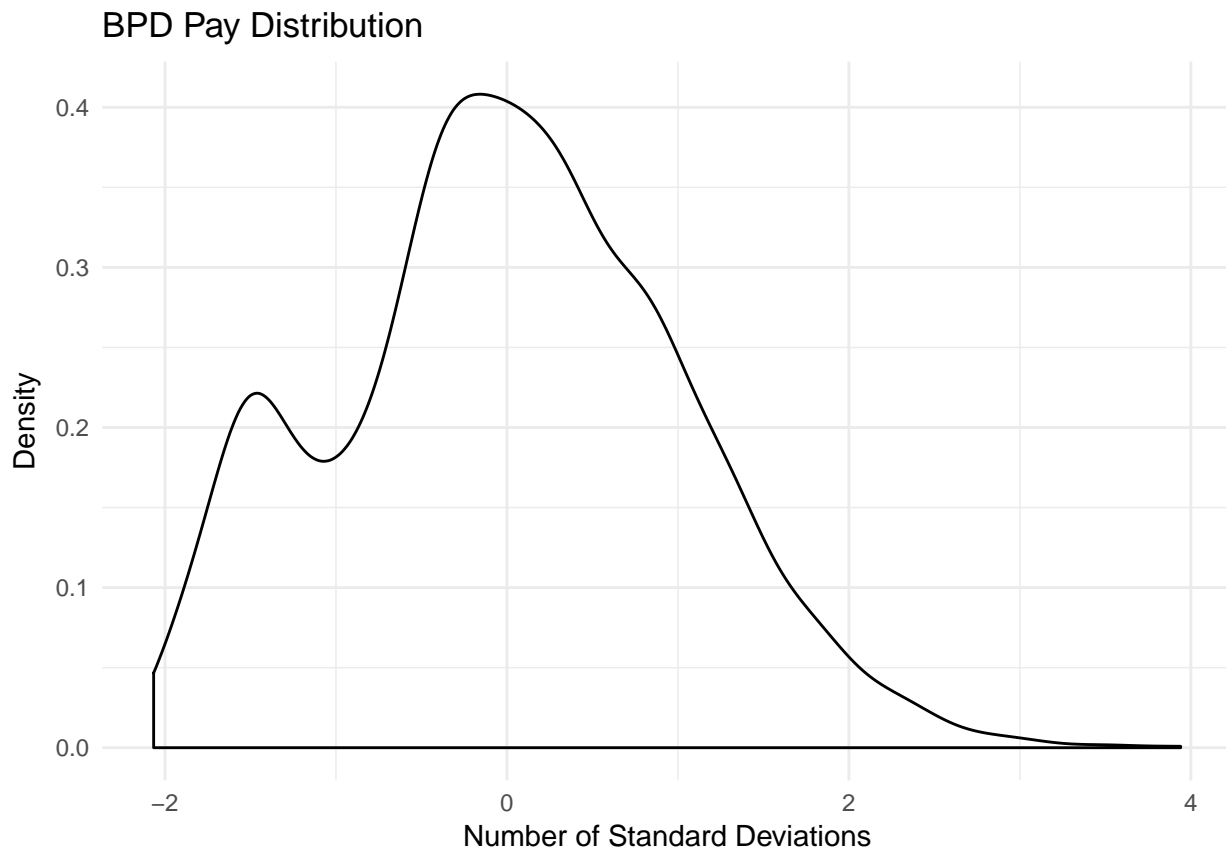
```
bpd_earners <- earnings_2017 %>% filter(`DEPARTMENT NAME` ==  
  "Boston Police Department")  
  
bpd_earners <- bpd_earners %>% mutate(num_std_devs = (.$`TOTAL EARNINGS` -  
  mean(bpd_earners$`TOTAL EARNINGS`))/sd(bpd_earners$`TOTAL EARNINGS`)) %>%  
  arrange(desc(.$`TOTAL EARNINGS`))  
  
head(bpd_earners)
```

```
## # A tibble: 6 x 9  
##   NAME `DEPARTMENT NAME` TITLE REGULAR OVERTIME EXTRA_PAY `TOTAL EARNINGS`  
##   <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl>  
## 1 Hose~ Boston Police D~ Poli~ 146894. 62696. 156642. 366233.  
## 2 Kerv~ Boston Police D~ Poli~ 125715. 66067. 150210. 341992.  
## 3 Lee,~ Boston Police D~ Poli~ 97414. 71669. 171093. 340176.  
## 4 Hass~ Boston Police D~ Poli~ 137104. 72158. 99536. 320224.  
## 5 McCo~ Boston Police D~ Poli~ 146894. 63708. 106072. 316674.  
## 6 Jose~ Boston Police D~ Poli~ 97414. 87746. 126997. 312156.  
## # ... with 2 more variables: INJURED <dbl>, num_std_devs <dbl>
```

```
bpd_2_sd <- bpd_earners %>% ggplot(aes(NAME, num_std_devs, color = (num_std_devs >  
  2))) + geom_point() + ggtitle("Officers who make more than $136071.46/year over the average") +  
  xlab("Officers") + ylab("Number of Standard Deviations") +  
  theme(axis.text.x = element_blank(), legend.title = element_blank())
```

```
bpd_3_sd <- bpd_earners %>% ggplot(aes(NAME, num_std_devs, color = (num_std_devs >  
  3))) + geom_point() + ggtitle("Officers who make more than $204107.19 / year over the average") +  
  xlab("Officers") + ylab("Number of Standard Deviations") +  
  theme(axis.text.x = element_blank(), legend.title = element_blank())
```

```
bpd_earners %>% filter(bpd_earners$`TOTAL EARNINGS` > 10000) %>%  
  ggplot(aes(.$num_std_devs)) + geom_density() + xlab("Number of Standard Deviations") +  
  ylab("Density") + ggtitle("BPD Pay Distribution") + theme_minimal()
```

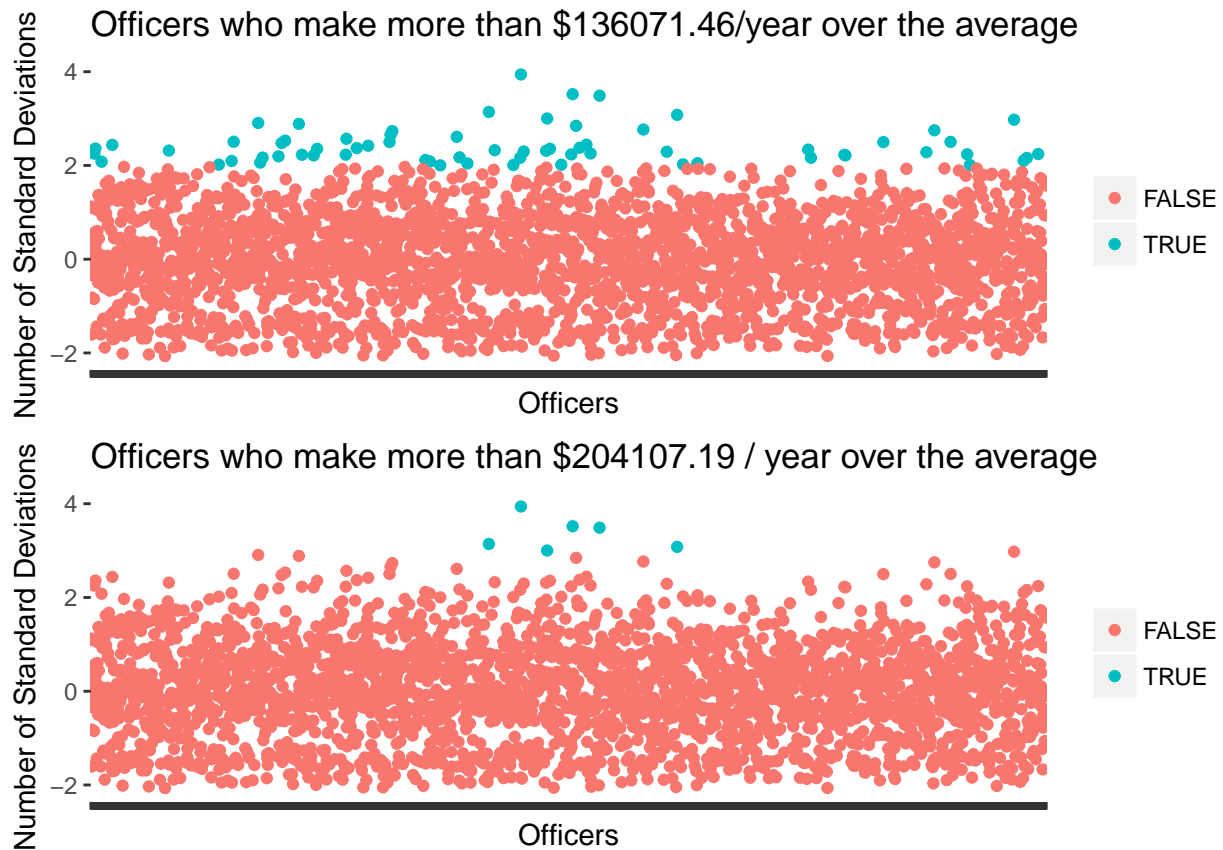



```
sd(bpd_earners$`TOTAL EARNINGS`)
```

```
## [1] 57576.17
```

We can see from the plot above that our Boston PD officers have pay that is roughly normally distributed around our mean. The graph above shows how many standard deviations officers are away from the mean. The mean itself is \$139345.4 per year and 1 standard deviation is \$57576.17 per year. Let's now look at our distribution again with an emphasis on our superearners:

```
gridExtra::grid.arrange(bpd_2_sd, bpd_3_sd)
```



We can see from above that there are a few people who earn exorbitant sums. What could cause this? The first thing that comes to mind is a high base salary.

```
bpd_earners %>% arrange(desc(.$REGULAR)) %>% head()
```

```
## # A tibble: 6 x 9
##   NAME `DEPARTMENT NAM~ TITLE REGULAR OVERTIME EXTRA_PAY `TOTAL EARNINGS`
##   <chr> <chr>           <chr>   <dbl>   <dbl>   <dbl>         <dbl>
## 1 Evan~ Boston Police D~ Comm~ 230000.     0    8846.     238846.
## 2 Gros~ Boston Police D~ Supn~ 199244.     0   26094.     225338.
## 3 Buck~ Boston Police D~ Supn~ 181983.     0   57978.     239961.
## 4 Manc~ Boston Police D~ Supn~ 181983.     0   57978.     239961.
## 5 Holm~ Boston Police D~ Supn~ 181983.     0   47087.     229069.
## 6 Ridg~ Boston Police D~ Supn~ 180369.     0   57726.     238094.
## # ... with 2 more variables: INJURED <dbl>, num_std_devs <dbl>
```

We can see from the table above that the highest paid person in the Boston Police Department is the Commissioner with a salary of about \$230,000 and a total earnings of about \$240,000. Even with this high base salary, he is only 1.68 standard deviations above the average officer pay of \$124254.60. In fact, the entire table of top base pay people have standard deviations less than 2, so we can see that base pay isn't what contributes to such high pay.

The next thing that comes to mind is a lot of overtime mixed with extra pay. This extra pay includes things like road detail and testifying in court. We could also check an assumption that those with a higher paygrade (e.g. captains and lieutenants) are likely to get more overtime money since their time and a half is generally a lot higher.

```
# create plot that shows overtime paid by base wage
bpd_overtime <- bpd_earners %>% filter(.$REGULAR > 20000) %>%
```

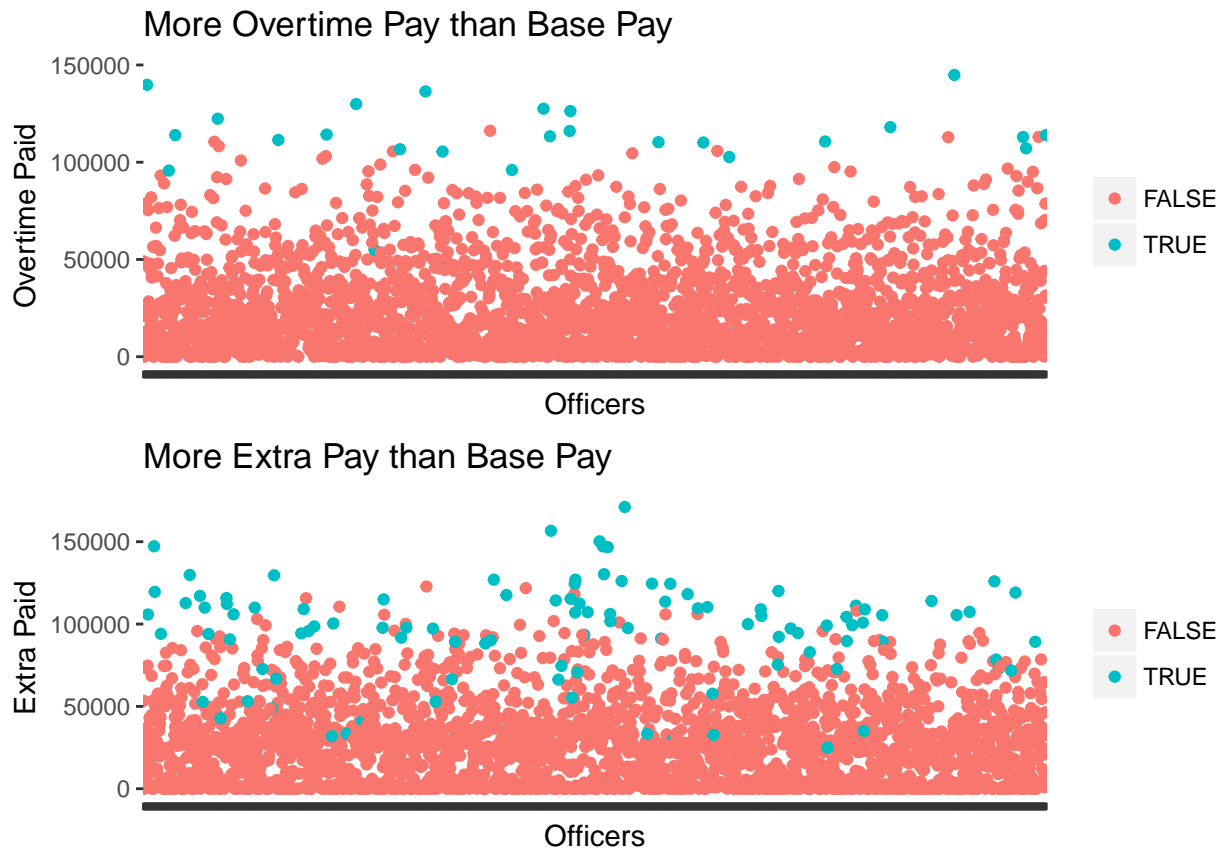
```

ggplot(aes(. $NAME, . $OVERTIME, color = (. $OVERTIME > . $REGULAR))) +
geom_point() + ggtitle("More Overtime Pay than Base Pay") +
theme(axis.text.x = element_blank(), legend.title = element_blank()) +
xlab("Officers") + ylab("Overtime Paid")

bpd_extra <- bpd_earners %>% filter(. $REGULAR > 20000) %>% ggplot(aes(. $NAME,
. $EXTRA_PAY, color = (. $EXTRA_PAY > . $REGULAR))) + geom_point() +
ggtitle("More Extra Pay than Base Pay") + theme(axis.text.x = element_blank(),
legend.title = element_blank()) + xlab("Officers") + ylab("Extra Paid")

gridExtra::grid.arrange(bpd_overtime, bpd_extra)

```



From the plots above we can see that there are quite a few people with more of their annual salary coming from overtime or extra pay. This is quite surprising.

If we can assume time and a half, then there are officers working their regular hours and then much more. There are a few who make over \$100,000 per year extra through just overtime.

We can also see that there are even more officers who make more than their base salary in extra pay. This extra pay was defined as things like court appearances, detail, retrograde pay and the Quinn education incentive which gives a small salary bump for having a criminal justice degree. Why could this be? After looking into it, I came across some Boston Globe reports. Here is a quote from <https://www.bostonglobe.com/metro/2017/06/20/for-some-boston-police-officers-extra-money-comes-easy/oS47lc7OuYyVZbTPBv1zQL/story.html> :

Quote: "

In what critics call an extreme example of a systemic problem, Lee Waiman bolstered his wages thanks to police union contracts that require that officers who work detail shifts or testify in court be paid a minimum of four hours, even if the assignment lasts only 30 minutes.

Last year, Lee earned \$58,600 by working more than 1,100 hours of overtime, according to a Globe review of police payroll records. Records show that Lee did not work 674 of those hours — more than 16 40-hour weeks — yet received time-and-a-half pay.

Most of Lee's overtime pay stemmed from court appearances that typically lasted no more than an hour, the Globe found. He was also paid for 2,771 hours for detail shifts, including 861 unworked hours. That allowed him to make close to \$130,000, a sum that did not include his overtime pay.

"It's a generous system," said Sam Tyler, president of the Boston Municipal Research Bureau, a fiscal watchdog group. "You're paid for hours you don't work. It isn't a new issue, but it's one that really does need stricter focus and management to control those costs."

"

Clearly this is a known pattern and has been looked at before. For example, last year the BPD overtime alone hit \$66.9 million.

<https://www.bostonglobe.com/metro/2018/02/16/bpd-captain-was-city-top-earner/iI4G1pnC7MOUODxo0XR4aO/story.html>

Looking back even further, I came across this boston.com article from 2007:

http://archive.boston.com/news/local/articles/2007/08/23/3_police_lieutenants_are_cited_for_alleged_detail_abuses/

Quote: "

The internal audit of shifts worked in 2005 concluded that Lieutenants Haseeb Hosein, Timothy Kervin, and Ghassoub Frangie engaged in untruthful reporting of hours, performed details that conflicted with a scheduled tour of duty, and received details through unauthorized means. Hosein and Kervin were also cited with breaking the law, but officials did not provide details on the alleged infractions.

The investigators accused Hosein, a 19-year veteran, of 203 violations that include 80 counts of inaccurate reporting on a detail card, 16 counts of receiving details outside the system, 24 counts of accepting a detail scheduled during his regular patrol shifts, and one count of breaking the law and conduct unbecoming an officer.

Kervin, a 20-year veteran, was charged with 191 violations that include 68 counts of inaccurate reporting on a detail card, 46 counts of receiving details outside the system, six counts of accepting a detail scheduled during his regular patrol shifts, and one count each of breaking the law and conduct unbecoming of an officer.

Frangie, a 29-year veteran, was charged with 84 violations that include 34 counts of inaccurate reporting on a detail card, 10 counts of accepting a detail scheduled during his regular patrol shifts, three counts of receiving details outside the system, and two counts of conduct unbecoming an officer.

"

```
# look at the top earning overtimers
bpd_earners %>% arrange(desc(.$OVERTIME)) %>% head(5)
```

```
## # A tibble: 5 x 9
##   NAME `DEPARTMENT NAM~ TITLE REGULAR OVERTIME EXTRA_PAY `TOTAL EARNINGS`
##   <chr> <chr>           <chr>    <dbl>    <dbl>    <dbl>          <dbl>
## 1 Sull~ Boston Police D~ Poli~ 110100.  144874.   28440.   283414.
## 2 Acos~ Boston Police D~ Poli~ 101681.  139806.   27517.   269004.
## 3 Fitz~ Boston Police D~ Poli~ 112121.  136404.   34853.   283378.
## 4 Deva~ Boston Police D~ Poli~ 116109.  129912.   28925.   274946.
## 5 Hold~ Boston Police D~ Poli~  95326.  127517.   32208.   255051.
## # ... with 2 more variables: INJURED <dbl>, num_std_devs <dbl>
```

look at the top earning extra pay people

```
bpd_earners %>% arrange(desc(.$EXTRA_PAY)) %>% head(5)
```

```
## # A tibble: 5 x 9
##   NAME `DEPARTMENT NAM~ TITLE REGULAR OVERTIME EXTRA_PAY `TOTAL EARNINGS`
##   <chr> <chr>           <chr>    <dbl>    <dbl>    <dbl>          <dbl>
## 1 Lee,~ Boston Police D~ Poli~  97414.   71669.  171093.   340176.
## 2 Hose~ Boston Police D~ Poli~ 146894.   62696.  156642.   366233.
## 3 Kerv~ Boston Police D~ Poli~ 125715.   66067.  150210.   341992.
## 4 Alme~ Boston Police D~ Poli~  86918.   24289.  147234.   259096.
## 5 King~ Boston Police D~ Poli~ 129531.   26444.  147170.   303145.
## # ... with 2 more variables: INJURED <dbl>, num_std_devs <dbl>
```

From the tables above we see some familiar names - including Lee Waiman from our first Boston Globe Article. We also see the names Haseeb Hosein and Timothy Kervin from our 2nd and 3rd articles. Clearly the city government is aware of the problem, but has not stopped it in at least 12 years.

That being said though, police are an important and vital part of society. Lets explore some of their valiant work.

Crime

```
# look at data set
```

```
crime %>% head()
```

```
## # A tibble: 6 x 17
```

```
##   INCIDENT_NUMBER OFFENSE_CODE OFFENSE_CODE_GRO~ OFFENSE_DESCRIP~ DISTRICT
##   <chr>           <chr>         <chr>         <chr>         <chr>
## 1 I182017604      03115      Investigate Pers~ INVESTIGATE PER~ B3
## 2 I182017601      00520      Residential Burg~ BURGLARY - RESI~ B2
## 3 I182017596      03831      Motor Vehicle Ac~ M/V - LEAVING S~ A1
## 4 I182017595      03802      Motor Vehicle Ac~ "M/V ACCIDENT -- <NA>
## 5 I182017594      01830      Drug Violation   DRUGS - SICK AS~ D14
## 6 I182017593      00361      Robbery          ROBBERY - OTHER  D4
## # ... with 12 more variables: REPORTING_AREA <int>, SHOOTING <chr>,
## #   OCCURRED_ON_DATE <dtm>, YEAR <int>, MONTH <int>, DAY_OF_WEEK <chr>,
## #   HOUR <int>, UCR_PART <chr>, STREET <chr>, Lat <dbl>, Long <dbl>,
## #   Location <chr>
```

```
# fix label
```

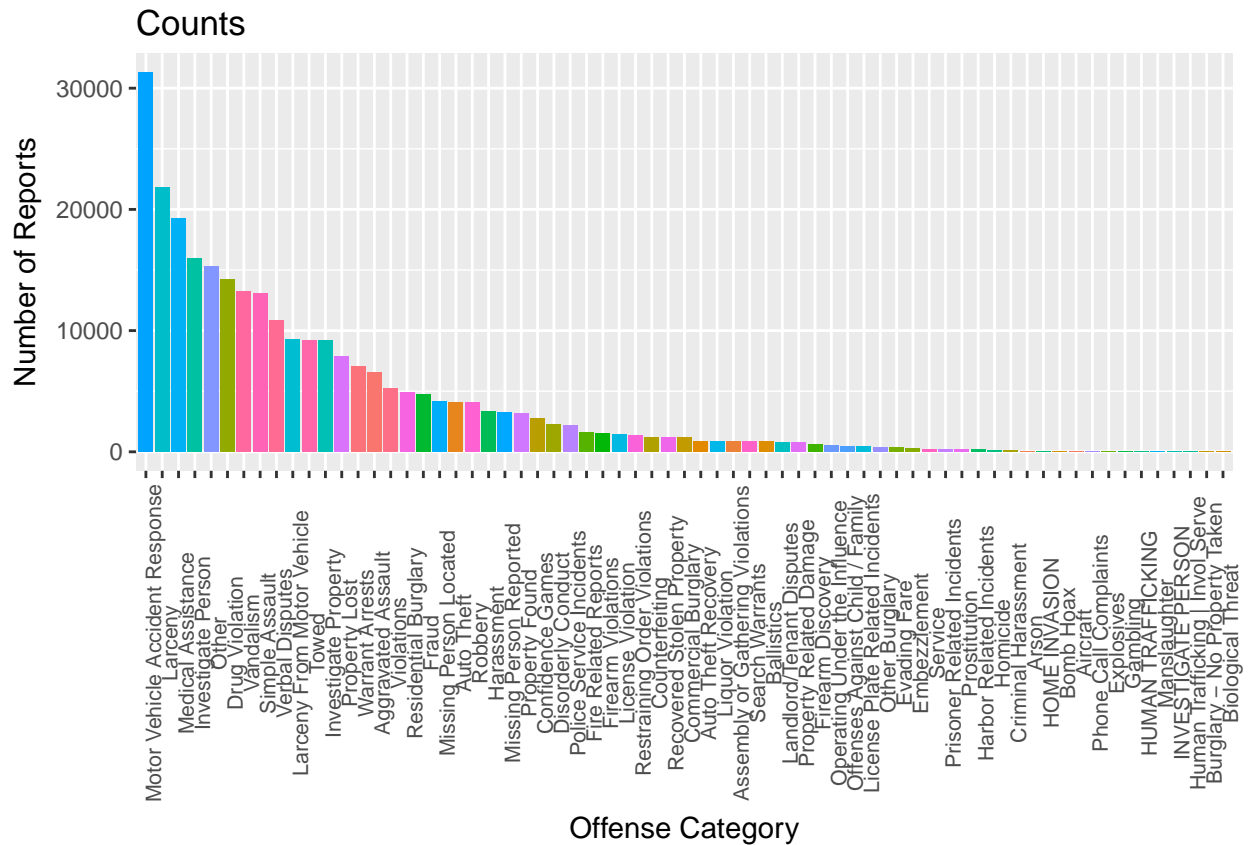
```
crime_by_count <- crime$OFFENSE_CODE_GROUP[crime$OFFENSE_CODE_GROUP ==
  "HUMAN TRAFFICKING - INVOLUNTARY SERVITUDE"] <- "Human Trafficking | Invol Serve"
```

```
# show number of each type of crime and proportions
```

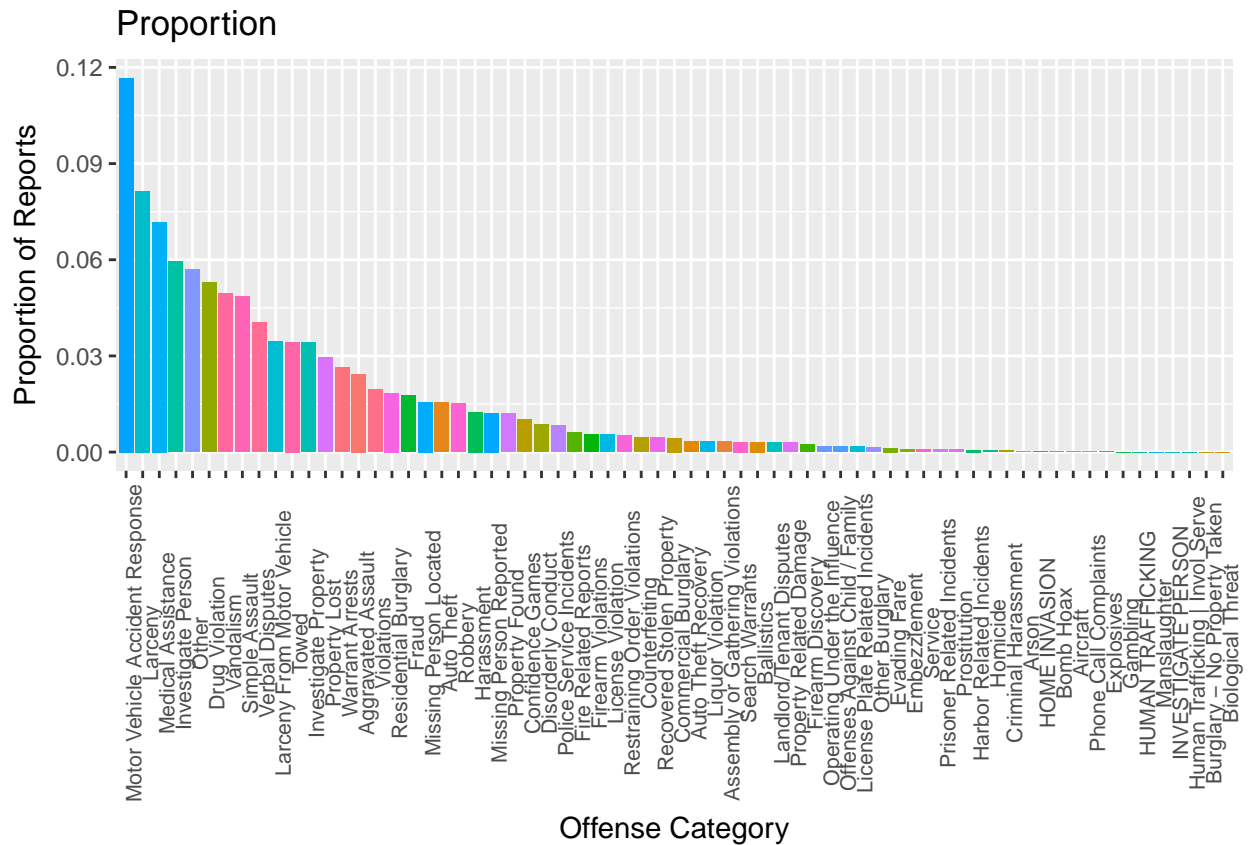
```
crime_by_count <- crime %>% group_by(.$OFFENSE_CODE_GROUP) %>%
  summarize(n = n(), prop = n/sum(length(crime$OFFENSE_CODE_GROUP)))
```

```
# plot
```

```
ggplot(crime_by_count, aes(fct_rev(fct_reorder(crime_by_count$`.$OFFENSE_CODE_GROUP`,
  crime_by_count$n)), crime_by_count$n)) + geom_bar(aes(fill = crime_by_count$`.$OFFENSE_CODE_GROUP`,
  stat = "identity")) + theme(axis.text.x = element_text(angle = 90,
  size = 8), legend.position = "none") + ylab("Number of Reports") +
  xlab("Offense Category") + ggtitle("Counts")
```



```
ggplot(crime_by_count, aes(fct_rev(fct_reorder(crime_by_count$`.$OFFENSE_CODE_GROUP`,
  crime_by_count$prop)), crime_by_count$prop)) + geom_bar(aes(fill = crime_by_count$`.$OFFENSE_CODE_G
  stat = "identity") + theme(axis.text.x = element_text(size = 8,
  angle = 90), legend.position = "none") + xlab("Offense Category") +
  ylab("Proportion of Reports") + ggtitle("Proportion")
```



As we see from above there are quite a few offense categories! Clearly Motor Vehicle Accident Response is the largest, followed by larceny (theft). This data set is quite interesting, so let's look at some more patterns.

Most Common Crimes by Time of Day

The tables below show the top 3 crimes that occur by hour. Understandably, the number 1 is consistently motor vehicle accident response, with consistently around ~10% of all crimes committed within that hour frame over 3 years of data (2015 - 2017). Other interesting bits:

- Simple assaults are more common during hours 1,2,3. This could be due to bars and nightlife.
- Vandalism commonly occurs during the hours of 4, 5, and 6 am. This makes sense, as vandals would be likely to strike at night.
- People got towed more often during the hours of 7, 8, and 9 am. This is likely cars that were left overnight.
- Drug Violations were most common during the hours of 16, 17, 18 and 19 (or 4,5,6,7 pm).

```
for (i in seq_along(1:23)) {  
  cat(sprintf("Hour: %d\n", i))  
  crime_hours <- crime %>% filter(crime$HOUR == i) %>% group_by(.$OFFENSE_CODE_GROUP)  
  
  crime_hours <- crime_hours %>% summarize(n = n(), prop = n/sum(length(crime_hours$OFFENSE_CODE_GROUP))  
    arrange(desc(.$n)) %>% head(3)  
  print(crime_hours)  
  cat("-----", "-----",  
      sep = "\n\n")  
}
```

```
## Hour: 1  
## # A tibble: 3 x 3  
##   `.$OFFENSE_CODE_GROUP`      n  prop  
##   <chr>                <int> <dbl>  
## 1 Motor Vehicle Accident Response  758 0.0962  
## 2 Simple Assault                657 0.0834  
## 3 Medical Assistance             600 0.0762  
## -----  
##  
## -----  
## Hour: 2  
## # A tibble: 3 x 3  
##   `.$OFFENSE_CODE_GROUP`      n  prop  
##   <chr>                <int> <dbl>  
## 1 Motor Vehicle Accident Response  831 0.128  
## 2 Simple Assault                604 0.0931  
## 3 Medical Assistance             486 0.0749  
## -----  
##  
## -----  
## Hour: 3  
## # A tibble: 3 x 3  
##   `.$OFFENSE_CODE_GROUP`      n  prop  
##   <chr>                <int> <dbl>  
## 1 Motor Vehicle Accident Response  529 0.136  
## 2 Medical Assistance             336 0.0861  
## 3 Simple Assault                241 0.0618  
## -----  
##  
## -----  
## Hour: 4
```

```

## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response  340 0.120
## 2 Medical Assistance                284 0.0999
## 3 Vandalism                        182 0.0640
## -----
##
## -----
## Hour: 5
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response  437 0.158
## 2 Medical Assistance                299 0.108
## 3 Vandalism                        190 0.0688
## -----
##
## -----
## Hour: 6
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response  739 0.175
## 2 Medical Assistance                349 0.0826
## 3 Vandalism                        265 0.0627
## -----
##
## -----
## Hour: 7
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1226 0.164
## 2 Towed                          954 0.128
## 3 Medical Assistance              477 0.0639
## -----
##
## -----
## Hour: 8
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1566 0.143
## 2 Towed                          1177 0.107
## 3 Larceny                         746 0.0681
## -----
##
## -----
## Hour: 9
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1462 0.118

```

```

## 2 Towed                                1068 0.0859
## 3 Larceny                               874 0.0703
## -----
##
## -----
## Hour: 10
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1435 0.105
## 2 Larceny                     1219 0.0892
## 3 Medical Assistance            1078 0.0789
## -----
##
## -----
## Hour: 11
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1431 0.104
## 2 Larceny                     1275 0.0927
## 3 Medical Assistance            1067 0.0776
## -----
##
## -----
## Hour: 12
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Larceny                     1702 0.109
## 2 Motor Vehicle Accident Response 1504 0.0961
## 3 Medical Assistance            1116 0.0713
## -----
##
## -----
## Hour: 13
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1562 0.110
## 2 Larceny                     1449 0.102
## 3 Medical Assistance            1100 0.0777
## -----
##
## -----
## Hour: 14
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1635 0.114
## 2 Larceny                     1614 0.113
## 3 Medical Assistance            1073 0.0749
## -----
##
##

```

```

## -----
## Hour: 15
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1763 0.127
## 2 Larceny                      1542 0.111
## 3 Medical Assistance             1017 0.0730
## -----
##
## -----
## Hour: 16
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 2061 0.122
## 2 Larceny                      1666 0.0989
## 3 Drug Violation                1471 0.0873
## -----
##
## -----
## Hour: 17
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 2204 0.126
## 2 Drug Violation                1878 0.107
## 3 Larceny                      1648 0.0938
## -----
##
## -----
## Hour: 18
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1943 0.114
## 2 Drug Violation                1752 0.102
## 3 Larceny                      1554 0.0908
## -----
##
## -----
## Hour: 19
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop
##   <chr>                <int> <dbl>
## 1 Motor Vehicle Accident Response 1655 0.111
## 2 Larceny                      1297 0.0873
## 3 Drug Violation                1269 0.0854
## -----
##
## -----
## Hour: 20
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n    prop

```

```

##   <chr>                                <int>  <dbl>
## 1 Motor Vehicle Accident Response  1451 0.109
## 2 Larceny                          1078 0.0810
## 3 Medical Assistance                996 0.0749
## -----
##
## -----
## Hour: 21
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                      <int>  <dbl>
## 1 Motor Vehicle Accident Response  1341 0.113
## 2 Medical Assistance                977 0.0825
## 3 Larceny                          784 0.0662
## -----
##
## -----
## Hour: 22
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                      <int>  <dbl>
## 1 Motor Vehicle Accident Response  1342 0.123
## 2 Medical Assistance                860 0.0790
## 3 Vandalism                        768 0.0705
## -----
##
## -----
## Hour: 23
## # A tibble: 3 x 3
##   `.$OFFENSE_CODE_GROUP`      n  prop
##   <chr>                      <int>  <dbl>
## 1 Motor Vehicle Accident Response  1075 0.121
## 2 Vandalism                        657 0.0737
## 3 Medical Assistance                641 0.0719
## -----
##
## -----

```

Shooting

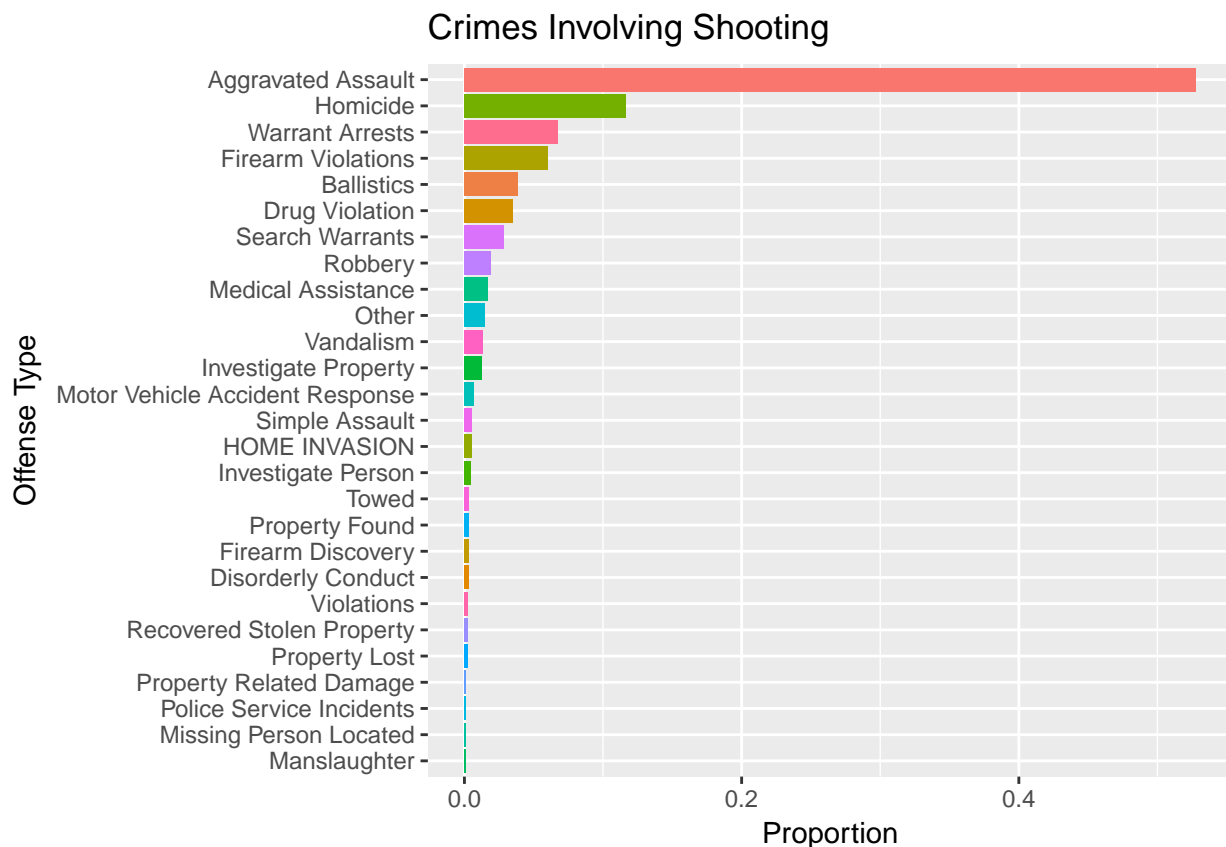
Lets take a look at the crimes involving shooting. The shooting column of the Crime Incidents Reports data indicates that a shooting took place. Lets take a look at the This data is from 2015 - 2018.

```
crime_shooting <- crime %>% filter(.$SHOOTING == "Y") %>% group_by(.$OFFENSE_CODE_GROUP)

crime_shooting <- crime_shooting %>% summarize(n = n(), prop = n/length(crime_shooting$OFFENSE_CODE_GROUP),
  arrange(desc(n))

crime_shooting <- crime_shooting %>% mutate(agg = cumsum(.$prop))

crime_shooting %>% ggplot() + geom_bar(aes(fct_reorder(crime_shooting$`.$OFFENSE_CODE_GROUP`,
  crime_shooting$n), crime_shooting$prop, fill = crime_shooting$`.$OFFENSE_CODE_GROUP`,
  stat = "identity")) + coord_flip() + theme(legend.position = "none") +
  xlab("Offense Type") + ylab("Proportion") + ggtitle("Crimes Involving Shooting")
```



As we can see from the table above, over 50% of crimes involving shooting were aggravated assaults. At around 10% of shooting crimes, there were 102 homicides. This is a pretty low number for 3 years in a major city - good job BPD!

Lets take a look at how the number of murders by shooting has changed over the years 2015 - 2018.

```
# filter crime by homicide, get counts and plot
murders <- crime %>% filter(.$OFFENSE_CODE_GROUP == "Homicide") %>%
  filter(.$SHOOTING == "Y") %>% filter(.$YEAR != 2018) %>%
  group_by(.$YEAR) %>% mutate(n = n()) %>% ggplot() + geom_line(aes(.$YEAR,
  n)) + theme_minimal() + xlab("Year") + ylab("Number of Crimes") +
```

```

ggtitle("Murder Counts") + scale_x_continuous(breaks = c(2014,
2015, 2016, 2017))

# filter crime by shooting, get counts and plot
overall_shooting <- crime %>% filter(.$SHOOTING == "Y") %>% filter(.$YEAR !=
2018) %>% group_by(.$YEAR) %>% mutate(n = n()) %>% ggplot() +
geom_line(aes(.$YEAR, n)) + theme_minimal() + xlab("Year") +
ylab("Number of Crimes") + ggtitle("Shooting Crimes") + scale_x_continuous(breaks = c(2014,
2015, 2016, 2017))

# filter crimes by aggravated assault with shooting, get
# counts and plot
agg_assaults <- crime %>% filter(.$OFFENSE_CODE_GROUP == "Aggravated Assault") %>%
filter(.$SHOOTING == "Y") %>% filter(.$YEAR != 2018) %>%
group_by(.$YEAR) %>% mutate(n = n()) %>% ggplot() + geom_line(aes(.$YEAR,
n)) + theme_minimal() + xlab("Year") + ylab("Number of Crimes") +
ggtitle("Aggravated Assaults") + scale_x_continuous(breaks = c(2014,
2015, 2016, 2017))

# look at gun recovery data
gun_recovery %>% head()

## # A tibble: 6 x 4
##   CollectionDate CrimeGunsRecovered GunsSurrenderedSafe~ BuybackGunsRecov~
##   <chr>                <int>                <int>                <int>
## 1 8/20/2014             2                  3                  1
## 2 8/21/2014             2                  0                  4
## 3 8/22/2014             0                  0                  2
## 4 8/25/2014             8                  3                  0
## 5 8/26/2014             9                  0                  0
## 6 8/27/2014             1                  0                  0

# add features for year and month to data frame
gun_recovery <- gun_recovery %>% mutate(Year = year(mdy(.$CollectionDate)),
Month = month(mdy(.$CollectionDate)))

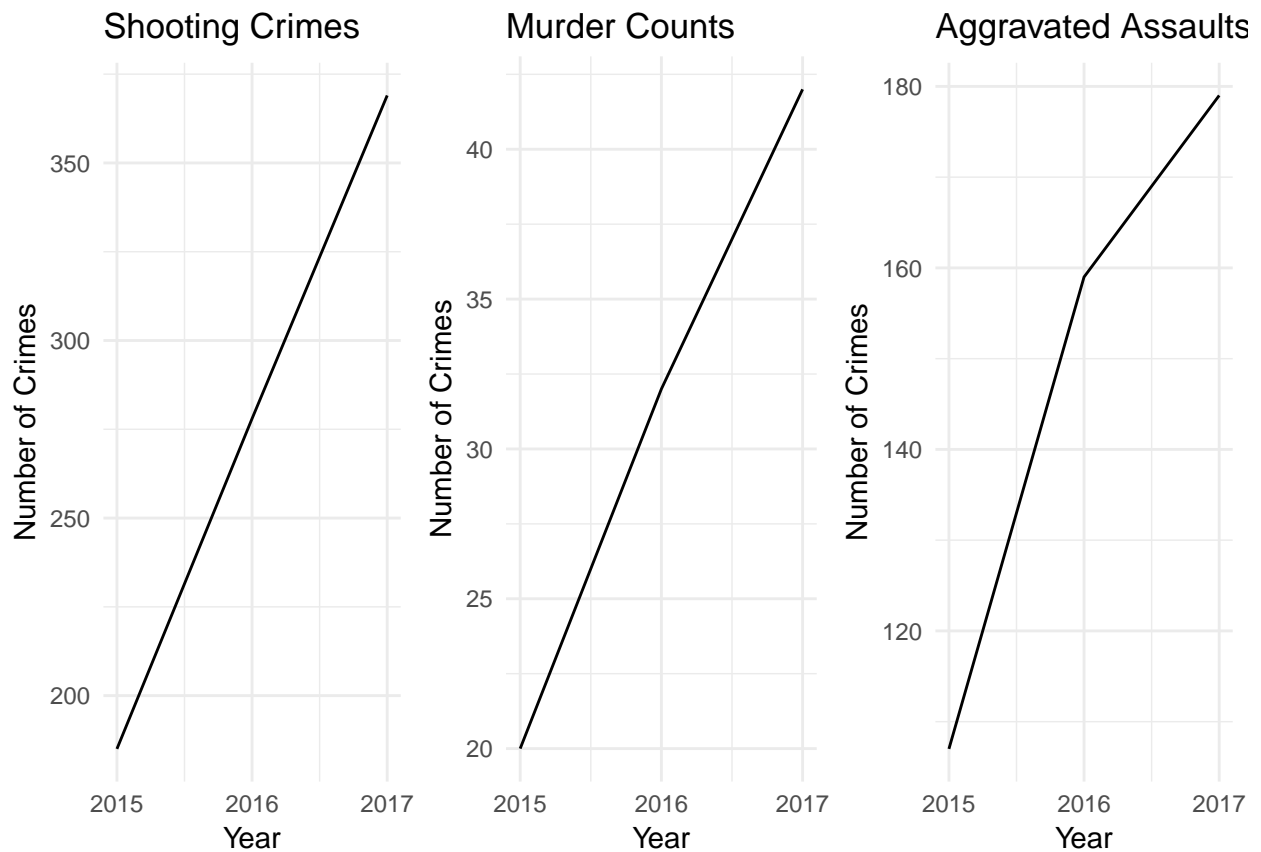
# create plot
gun_crimes <- gun_recovery %>% ggplot() + geom_jitter(aes(gun_recovery$Year,
gun_recovery$CrimeGunsRecovered)) + theme_minimal() + xlab("Year") +
ylab("Guns Recovered") + ggtitle("Guns Recovered from Crimes")

# plot voluntary surrender of guns
vol_surr <- gun_recovery %>% ggplot() + geom_jitter(aes(gun_recovery$Year,
gun_recovery$GunsSurrenderedSafeguarded)) + theme_minimal() +
xlab("Year") + ylab("Number of Guns") + ggtitle("Number of Guns Voluntarily Surrendered") +
scale_x_continuous(breaks = c(2014, 2015, 2016, 2017))

# plot results of gun buyback program
gun_buyback <- gun_recovery %>% ggplot() + geom_jitter(aes(gun_recovery$Year,
gun_recovery$BuybackGunsRecovered)) + theme_minimal() + xlab("Year") +
ylab("Number of Guns Bought") + ggtitle("Gun Buyback Program")

```

```
# plot
gridExtra::grid.arrange(overall_shooting, murders, agg_assaults,
  nrow = 1)
```

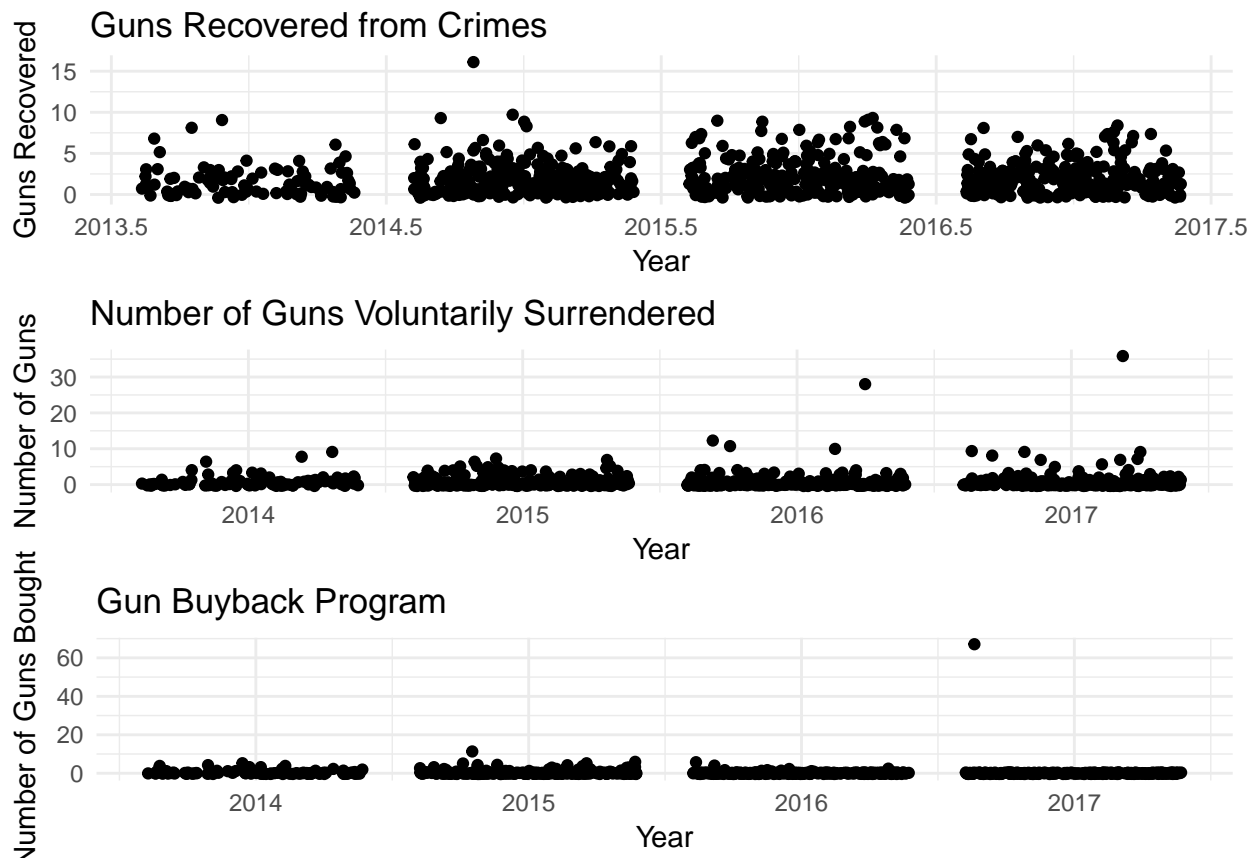


```
gridExtra::grid.arrange(gun_crimes, vol_surr, gun_buyback, nrow = 3)
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

From the first set of graphs we see the following:

- The number of shooting crimes jumps by about 80 each year above the previous year
- The number of murders jumps by about 10 each year more than the previous year
- The number of aggravated assaults rises, but the rate of increase has been slowed significantly between 2016 and 2017

From the second set of graphs related to gun recovery we see the following:

- The number of guns retrieved from crimes increased from 2014 to 2015, but there wasn't a large increase for the years 2015 - 2017
- The number of guns voluntarily surrendered seems to have been relatively consistent except for a few outliers (such as the 30+ guns surrendered in 2017)
- The gun buyback program has returned less guns than police work, but is still making a dent in the number of guns on the street. There is one remarkable data point in which 60+ guns were recovered.

Economic Indicators

Lets take a look at the Economic Indicators Dataset. This set contains information on

- Tourism/Flights
- Hotel Market
- Labor Market
- Real Estate: Board Approved Development Projects (Pipeline)
- Real Estate Market: Housing

```
econ_indicators %>% head()
```

```
## # A tibble: 6 x 19
##   Year Month logan_passengers logan_intl_flights hotel_occup_rate
##   <int> <int>          <int>          <int>          <dbl>
## 1  2013     1          2019662          2986          0.572
## 2  2013     2          1878731          2587          0.645
## 3  2013     3          2469155          3250          0.819
## 4  2013     4          2551246          3408          0.855
## 5  2013     5          2676291          3240          0.858
## 6  2013     6          2824862          3402          0.911
## # ... with 14 more variables: hotel_avg_daily_rate <dbl>,
## #   total_jobs <int>, unemp_rate <dbl>, labor_force_part_rate <dbl>,
## #   pipeline_unit <int>, pipeline_total_dev_cost <dbl>,
## #   pipeline_sqft <int>, pipeline_const_jobs <dbl>, foreclosure_pet <int>,
## #   foreclosure_deeds <int>, med_housing_price <int>,
## #   housing_sales_vol <int>, new_housing_const_permits <int>,
## #   `new-affordable_housing_permits` <int>
```

```
names(econ_indicators)
```

```
## [1] "Year"                "Month"
## [3] "logan_passengers"    "logan_intl_flights"
## [5] "hotel_occup_rate"    "hotel_avg_daily_rate"
## [7] "total_jobs"          "unemp_rate"
## [9] "labor_force_part_rate" "pipeline_unit"
## [11] "pipeline_total_dev_cost" "pipeline_sqft"
## [13] "pipeline_const_jobs"   "foreclosure_pet"
## [15] "foreclosure_deeds"     "med_housing_price"
## [17] "housing_sales_vol"     "new_housing_const_permits"
## [19] "new-affordable_housing_permits"
```

```
unique(econ_indicators$total_jobs)
```

```
## [1] 0
```

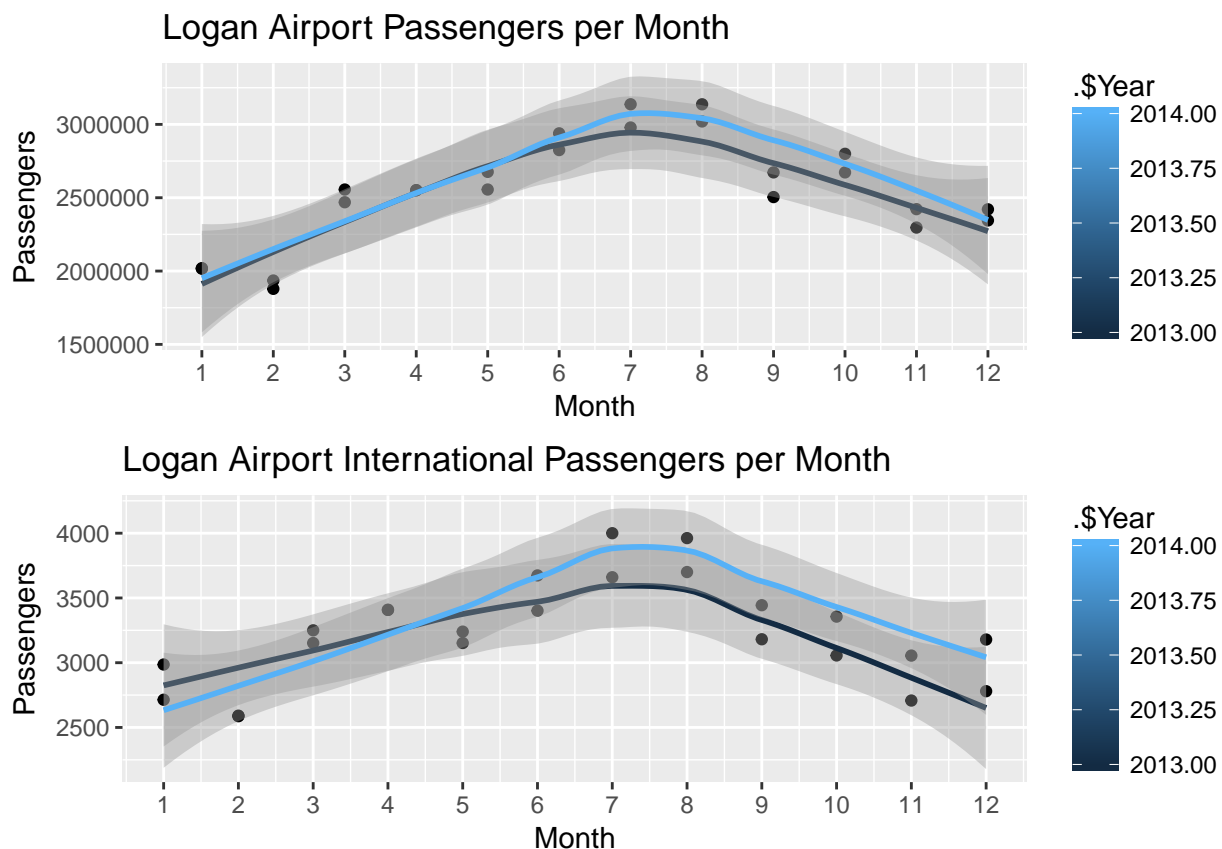
Flights / Tourism

```
logan_pass <- econ_indicators %>% filter(.$logan_passengers >
  0) %>% ggplot() + geom_point(aes(.$Month, .$logan_passengers)) +
  geom_smooth(aes(.$Month, .$logan_passengers, group = .$Year,
    color = .$Year)) + xlab("Month") + ylab("Passengers") +
  ggtitle("Logan Airport Passengers per Month") + scale_x_continuous(breaks = c(1,
    2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

logan_intl_pass <- econ_indicators %>% filter(.$logan_intl_flights >
  0) %>% ggplot() + geom_point(aes(.$Month, .$logan_intl_flights)) +
  geom_smooth(aes(.$Month, .$logan_intl_flights, group = .$Year,
    color = .$Year)) + xlab("Month") + ylab("Passengers") +
  ggtitle("Logan Airport International Passengers per Month") +
  scale_x_continuous(breaks = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
    10, 11, 12))

gridExtra::grid.arrange(logan_pass, logan_intl_pass)
```

```
## `geom_smooth()` using method = 'loess'
## `geom_smooth()` using method = 'loess'
```



From the graphs above we see the following trends:

- 2014 was a slightly better year for tourism. This makes sense from an economic perspective as we will see in the graphs that follow.
- Flights jump about 50% from the beginning of the year to the summer and then level off as it gets

colder

- There is roughly a thousand times as many domestic passengers as there are international passengers

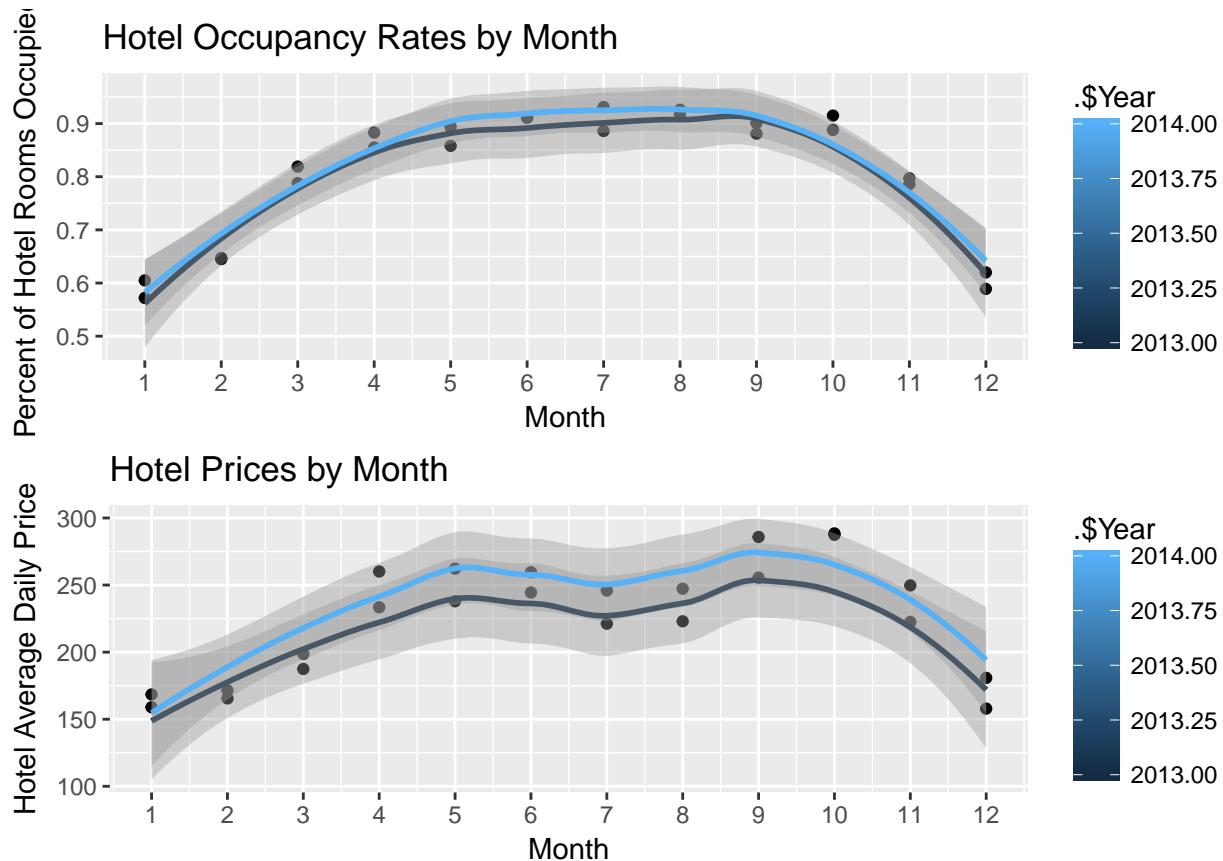
Hotels

```
hotel_occ <- econ_indicators %>% filter(.$hotel_occup_rate >
0) %>% ggplot() + geom_point(aes(.$Month, .$hotel_occup_rate)) +
geom_smooth(aes(.$Month, .$hotel_occup_rate, group = .$Year,
color = .$Year)) + xlab("Month") + ylab("Percent of Hotel Rooms Occupied") +
ggtitle("Hotel Occupancy Rates by Month") + scale_x_continuous(breaks = c(1,
2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

hotel_prices <- econ_indicators %>% filter(.$hotel_avg_daily_rate >
0) %>% ggplot() + geom_point(aes(.$Month, .$hotel_avg_daily_rate)) +
geom_smooth(aes(.$Month, .$hotel_avg_daily_rate, group = .$Year,
color = .$Year)) + xlab("Month") + ylab("Hotel Average Daily Price") +
ggtitle("Hotel Prices by Month") + scale_x_continuous(breaks = c(1,
2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

gridExtra::grid.arrange(hotel_occ, hotel_prices)
```

```
## `geom_smooth()` using method = 'loess'
## `geom_smooth()` using method = 'loess'
```



From the graphs above we see the following trends:

- 2014 had similar occupation rates to 2013
- The average price of hotels increased from 2013 to 2014
- The hotels are around half full during the winter and almost completely full during the summer
- Christmas is not that popular of a time to book a hotel in Boston

- The price of hotels increases as occupation increases

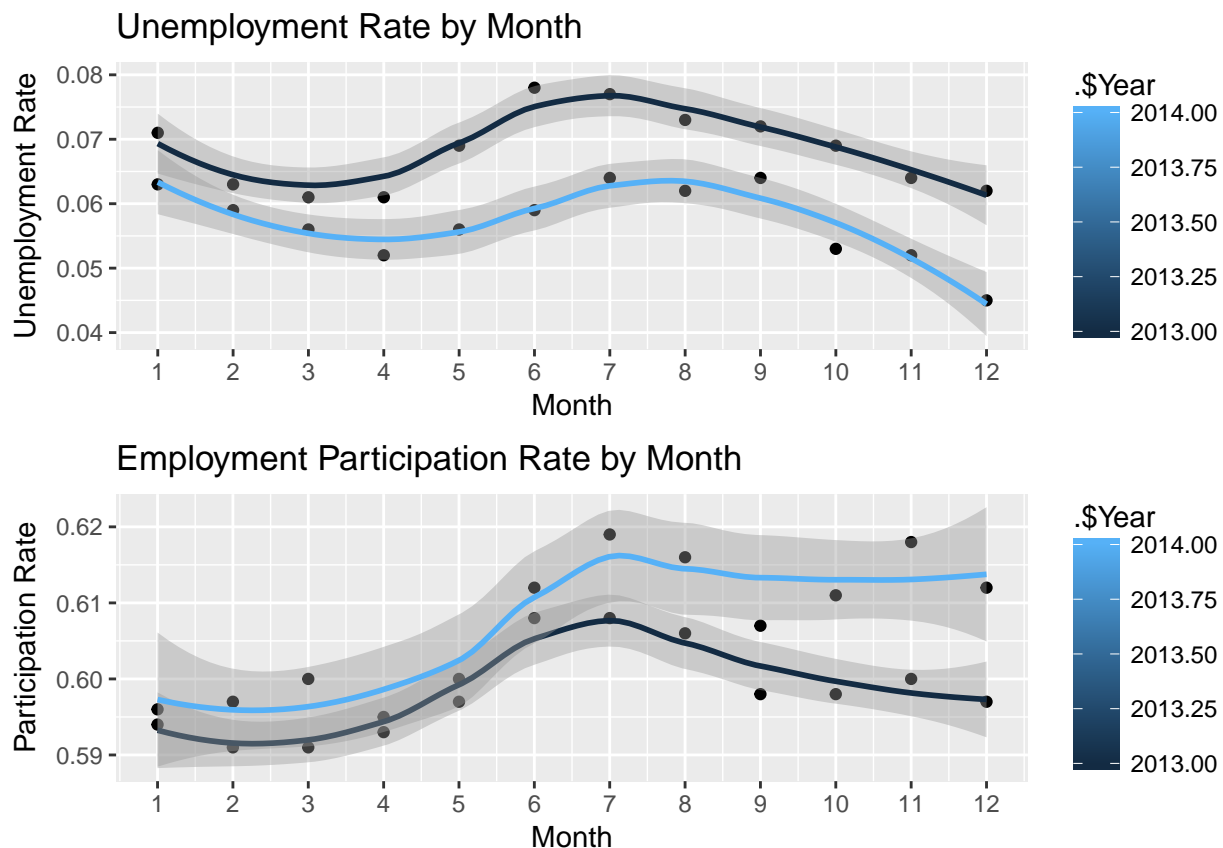
Job Market

```
unemp_rate <- econ_indicators %>% filter(.$unemp_rate > 0.001) %>%
  ggplot() + geom_point(aes(.$Month, .$unemp_rate)) + geom_smooth(aes(.$Month,
    .$unemp_rate, group = .$Year, color = .$Year)) + xlab("Month") +
  ylab("Unemployment Rate") + ggtitle("Unemployment Rate by Month") +
  scale_x_continuous(breaks = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
    10, 11, 12))

part_rate <- econ_indicators %>% filter(.$labor_force_part_rate >
  0.1) %>% ggplot() + geom_point(aes(.$Month, .$labor_force_part_rate)) +
  geom_smooth(aes(.$Month, .$labor_force_part_rate, group = .$Year,
    color = .$Year)) + xlab("Month") + ylab("Participation Rate") +
  ggtitle("Employment Participation Rate by Month") + scale_x_continuous(breaks = c(1,
    2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

gridExtra::grid.arrange(unemp_rate, part_rate)
```

```
## `geom_smooth()` using method = 'loess'
## `geom_smooth()` using method = 'loess'
```



The graphs above show the following trends:

- 2013 had a higher unemployment rate than 2014
- Unemployment dropped on average 1%-2% between the years, dependent on the month
- Employees were more likely to participate in the employment market in 2014 than 2013.
- Unemployment rates are generally around 6-8%

- Market participation is roughly 60%

Real Estate Board Approved Development Projects

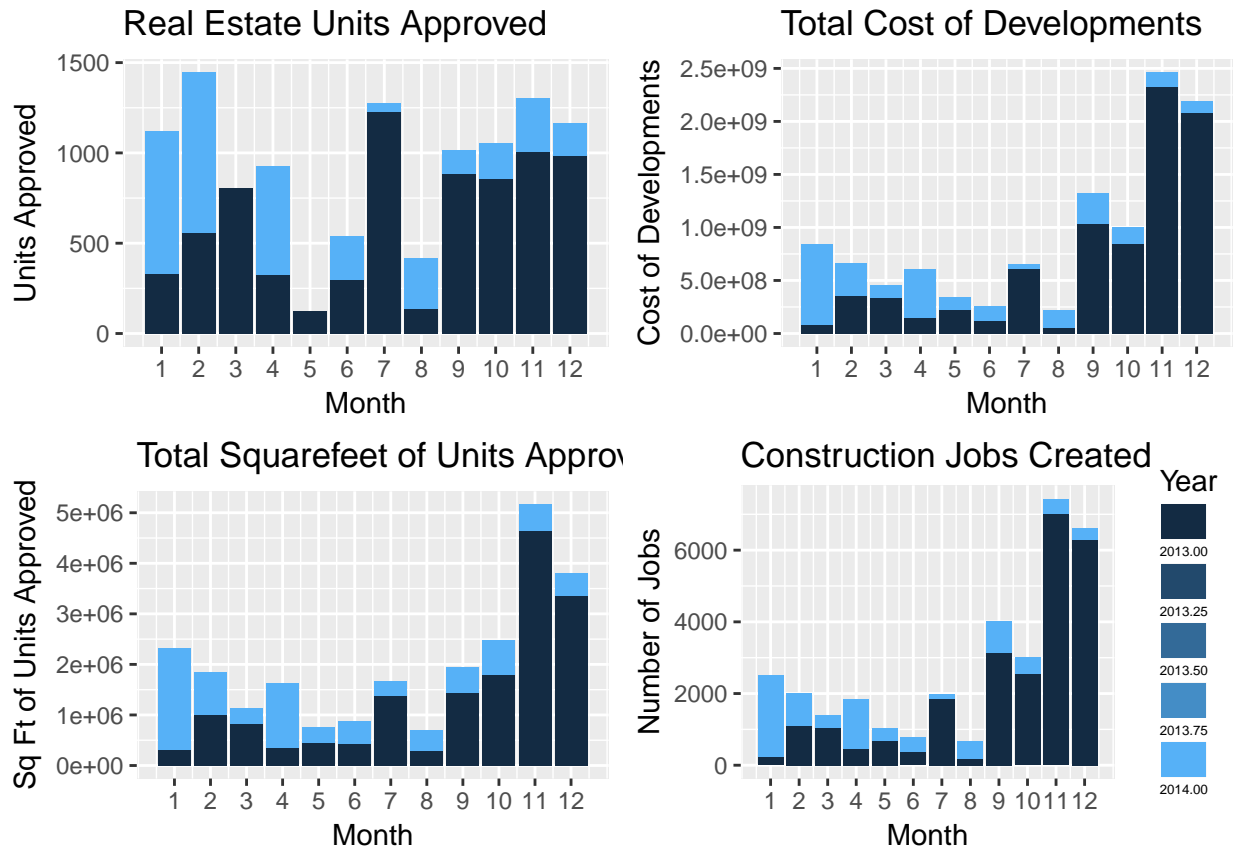
```
units_approved <- econ_indicators %>% filter(. $pipeline_unit >
0) %>% group_by(. $Year) %>% ggplot() + geom_bar(aes(. $Month,
. $pipeline_unit, fill = . $Year), position = "stack", stat = "identity") +
xlab("Month") + ylab("Units Approved") + ggtitle("Real Estate Units Approved") +
theme(legend.position = "none") + scale_x_continuous(breaks = c(1,
2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

cost_of_devs <- econ_indicators %>% filter(. $pipeline_total_dev_cost >
0) %>% group_by(. $Year) %>% ggplot() + geom_bar(aes(. $Month,
. $pipeline_total_dev_cost, fill = . $Year), position = "stack",
stat = "identity") + xlab("Month") + ylab("Cost of Developments") +
ggtitle("Total Cost of Developments") + theme(legend.position = "none") +
scale_x_continuous(breaks = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
10, 11, 12))

tot_sqft_devs <- econ_indicators %>% filter(. $pipeline_sqft >
0) %>% group_by(. $Year) %>% ggplot() + geom_bar(aes(. $Month,
. $pipeline_sqft, fill = . $Year), position = "stack", stat = "identity") +
xlab("Month") + ylab("Sq Ft of Units Approved") + ggtitle("Total Squarefeet of Units Approved") +
theme(legend.position = "none") + scale_x_continuous(breaks = c(1,
2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

const_jobs <- econ_indicators %>% filter(. $pipeline_const_jobs >
0) %>% group_by(. $Year) %>% ggplot() + geom_bar(aes(. $Month,
. $pipeline_const_jobs, fill = . $Year), position = "stack",
stat = "identity") + xlab("Month") + ylab("Number of Jobs") +
ggtitle("Construction Jobs Created") + guides(fill = guide_legend(keywidth = 1,
keyheight = 1, title = "Year", label.position = "bottom")) +
theme(legend.text = element_text(size = 5)) + scale_x_continuous(breaks = c(1,
2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

gridExtra::grid.arrange(units_approved, cost_of_devs, tot_sqft_devs,
const_jobs)
```



The graphs above show the following trends:

- 2014 had a significant amount more development than 2013
- The boom in development also lead to an increase in construction jobs
- The units that were developed also took up more space

Housing

```
forc_pets <- econ_indicators %>% filter(.$foreclosure_pet > 0) %>%
  group_by(.$Year) %>% ggplot() + geom_bar(aes(.$Month, .$foreclosure_pet,
  fill = .$Year), position = "stack", stat = "identity") +
  xlab("Month") + ylab("Number of Foreclosure Petitions") +
  ggtitle("Foreclosure Petitions") + theme(legend.position = "none") +
  scale_x_continuous(breaks = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
    10, 11, 12))

forc_deeds <- econ_indicators %>% filter(.$foreclosure_deeds >
  0) %>% group_by(.$Year) %>% ggplot() + geom_bar(aes(.$Month,
  .$foreclosure_deeds, fill = .$Year), position = "stack",
  stat = "identity") + xlab("Month") + ylab("Number of Foreclosure Deeds") +
  ggtitle("Foreclosure Deeds") + guides(fill = guide_legend(keywidth = 1,
  keyheight = 1, title = "Year", label.position = "bottom")) +
  theme(legend.text = element_text(size = 5)) + scale_x_continuous(breaks = c(1,
  2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

housing_prices <- econ_indicators %>% filter(.$med_housing_price >
  0) %>% group_by(.$Year) %>% ggplot() + geom_smooth(aes(.$Month,
  .$med_housing_price, color = .$Year, group = .$Year), se = FALSE) +
  xlab("Month") + ylab("Median Housing Price") + ggtitle("Median Housing Sales Price") +
  guides(fill = guide_legend(keywidth = 1, keyheight = 1, title = "Year",
    label.position = "bottom")) + theme(legend.text = element_text(size = 5)) +
  scale_x_continuous(breaks = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
    10, 11, 12))

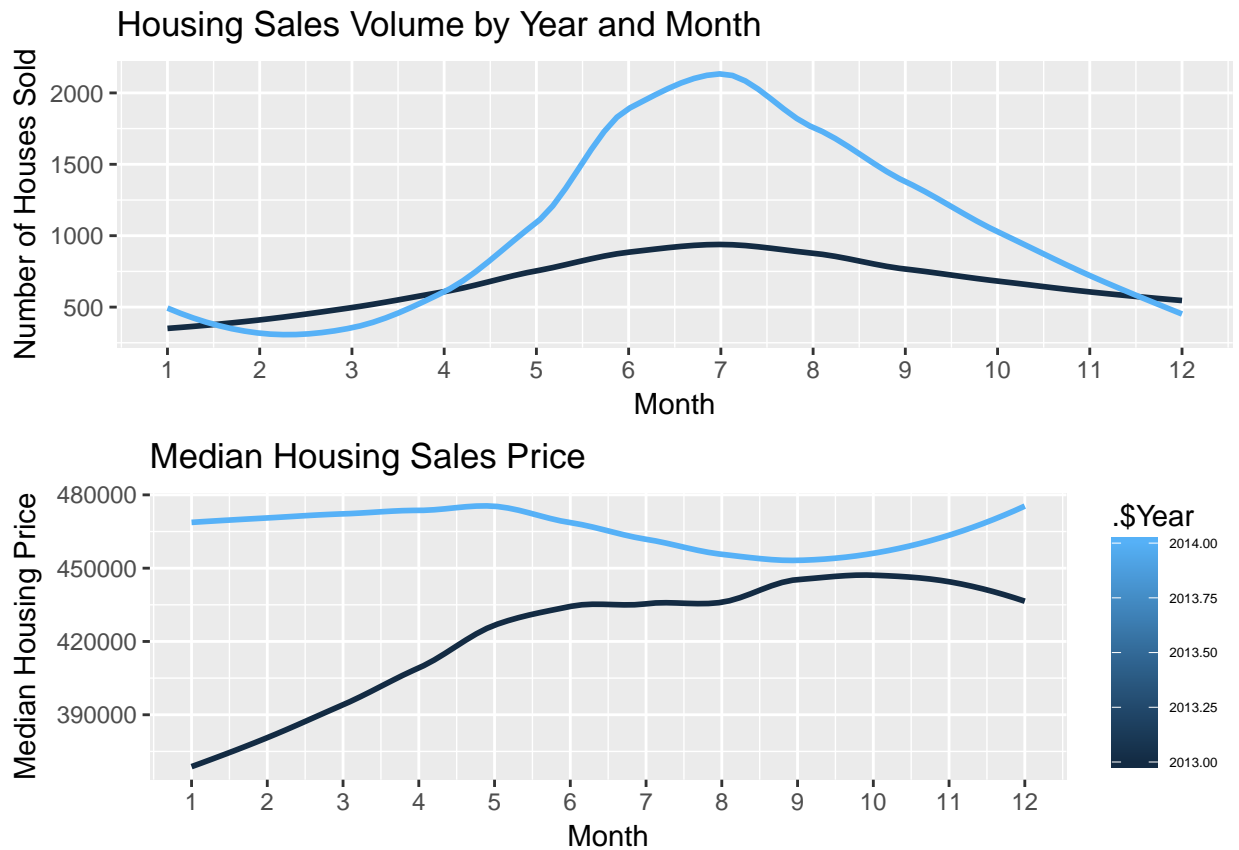
housing_sales <- econ_indicators %>% filter(.$housing_sales_vol >
  0) %>% group_by(.$Year) %>% ggplot() + geom_smooth(aes(.$Month,
  .$housing_sales_vol, color = .$Year, group = .$Year), se = FALSE) +
  xlab("Month") + ylab("Number of Houses Sold") + ggtitle("Housing Sales Volume by Year and Month") +
  theme(legend.position = "none") + scale_x_continuous(breaks = c(1,
  2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))

housing_const_permits <- econ_indicators %>% filter(.$new_housing_const_permits >
  0) %>% group_by(.$Year) %>% ggplot() + geom_bar(aes(.$Month,
  .$new_housing_const_permits, fill = .$Year), position = "stack",
  stat = "identity") + xlab("Month") + ylab("Number of Housing Construction Permits") +
  ggtitle("Housing Construction") + theme(legend.position = "none") +
  scale_x_continuous(breaks = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
    10, 11, 12))

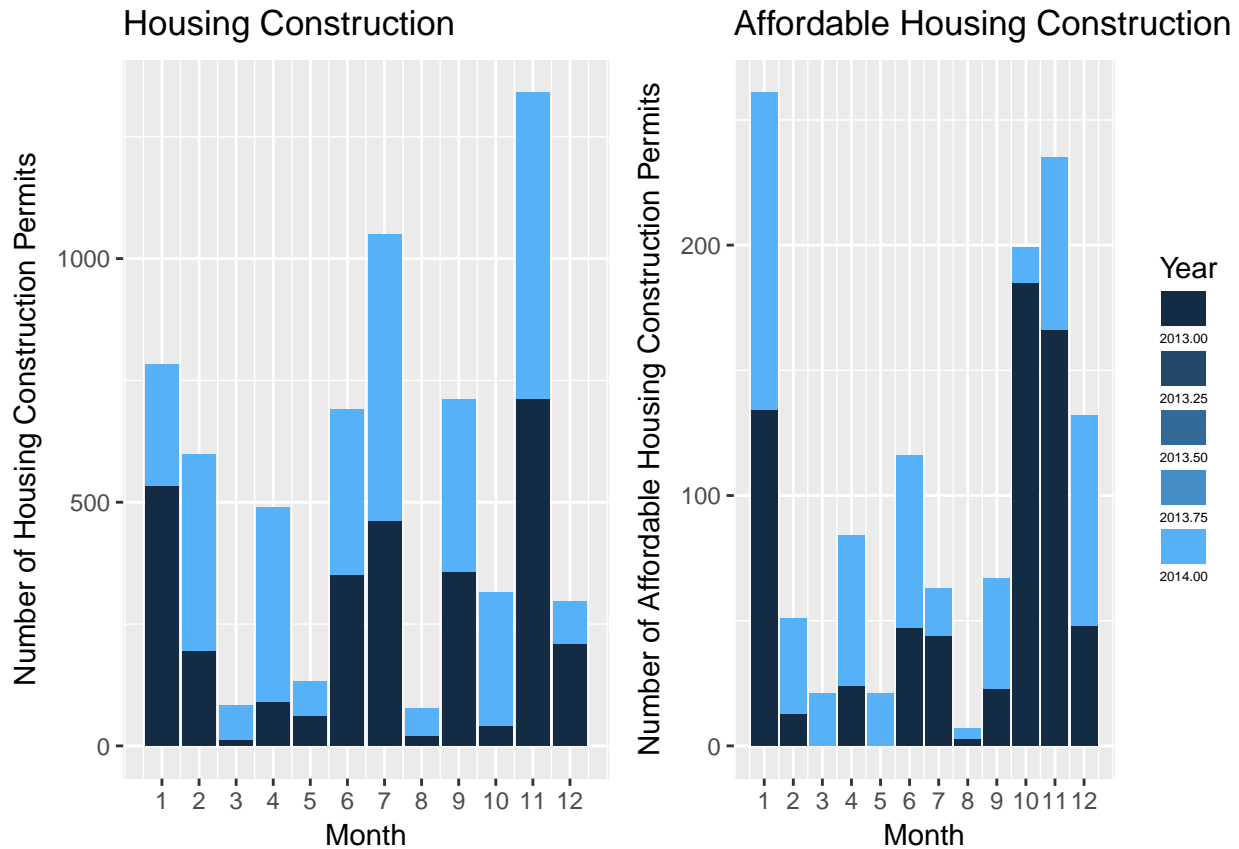
affordable_permits <- econ_indicators %>% filter(.$`new-affordable_housing_permits` >
  0) %>% group_by(.$Year) %>% ggplot() + geom_bar(aes(.$Month,
  .$`new-affordable_housing_permits`, fill = .$Year), position = "stack",
  stat = "identity") + xlab("Month") + ylab("Number of Affordable Housing Construction Permits") +
  ggtitle("Affordable Housing Construction") + guides(fill = guide_legend(keywidth = 1,
  keyheight = 1, title = "Year", label.position = "bottom")) +
  theme(legend.text = element_text(size = 5)) + scale_x_continuous(breaks = c(1,
  2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12))
```

```
gridExtra::grid.arrange(housing_sales, housing_prices, nrow = 2)
```

```
## `geom_smooth()` using method = 'loess'
## `geom_smooth()` using method = 'loess'
```



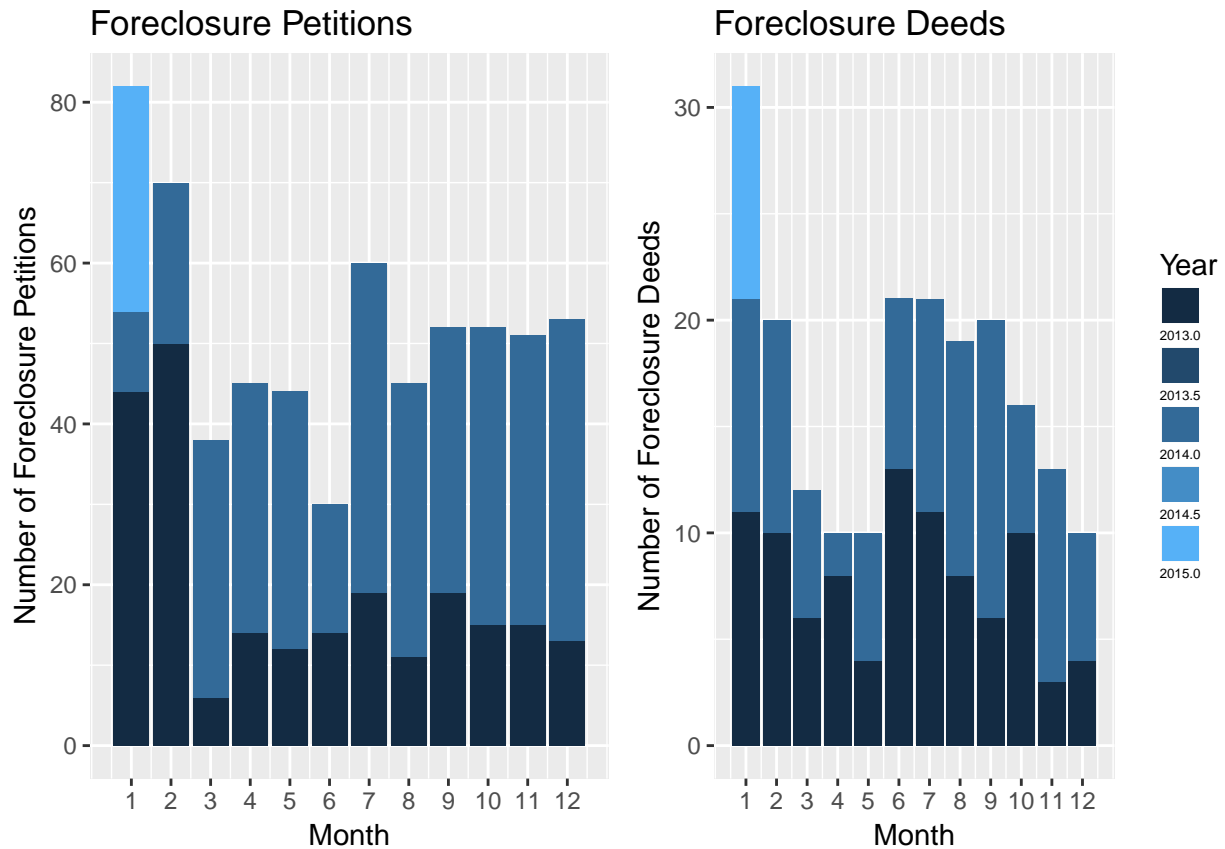
```
gridExtra::grid.arrange(housing_const_permits, affordable_permits,
  nrow = 1)
```



From the above graphs we can see the following trends:

- Housing sales increase during the summer
- The housing market had over twice the number of sales in 2014 and it did in 2013
- Houses sold in 2013 started off quite low at the beginning of the year and gradually increased in price
- Housing sales prices were roughly uniform over 2014
- Housing construction (affordable and regular) increased a good amount in 2014
- The number of construction permits for regular housing was quite high in November
- There was a lot more affordable housing being built in 2014 than in 2013

```
gridExtra::grid.arrange(forc_pets, forc_deeds, nrow = 1)
```



A foreclosure petition is when a lender (typically a bank) sues a delinquent tenant by filing a court document for foreclosure. This petition for foreclosure is then delivered to the homeowner along with a court summons.

A foreclosure deed is when a lender accepts the deed (document stating ownership) of a property instead of foreclosing on a house.

The graphs above show the following trends:

- The amount of foreclosures rose in 2014 from 2013
- Foreclosure notices in 2013 and 2014 were generally distributed between all the months of the year
- Foreclosure notices in 2015 were generally given in January
- 2015 seems to have had less foreclosures than 2013 and 2014 (perhaps a sign of the housing market crash of 08 recovering)