

```
#first, I open the file and make a simple plot of all crimes by district
setwd("/Users/igor/MOOCs/data_manipulation/datasci_course_materials/assignment6")

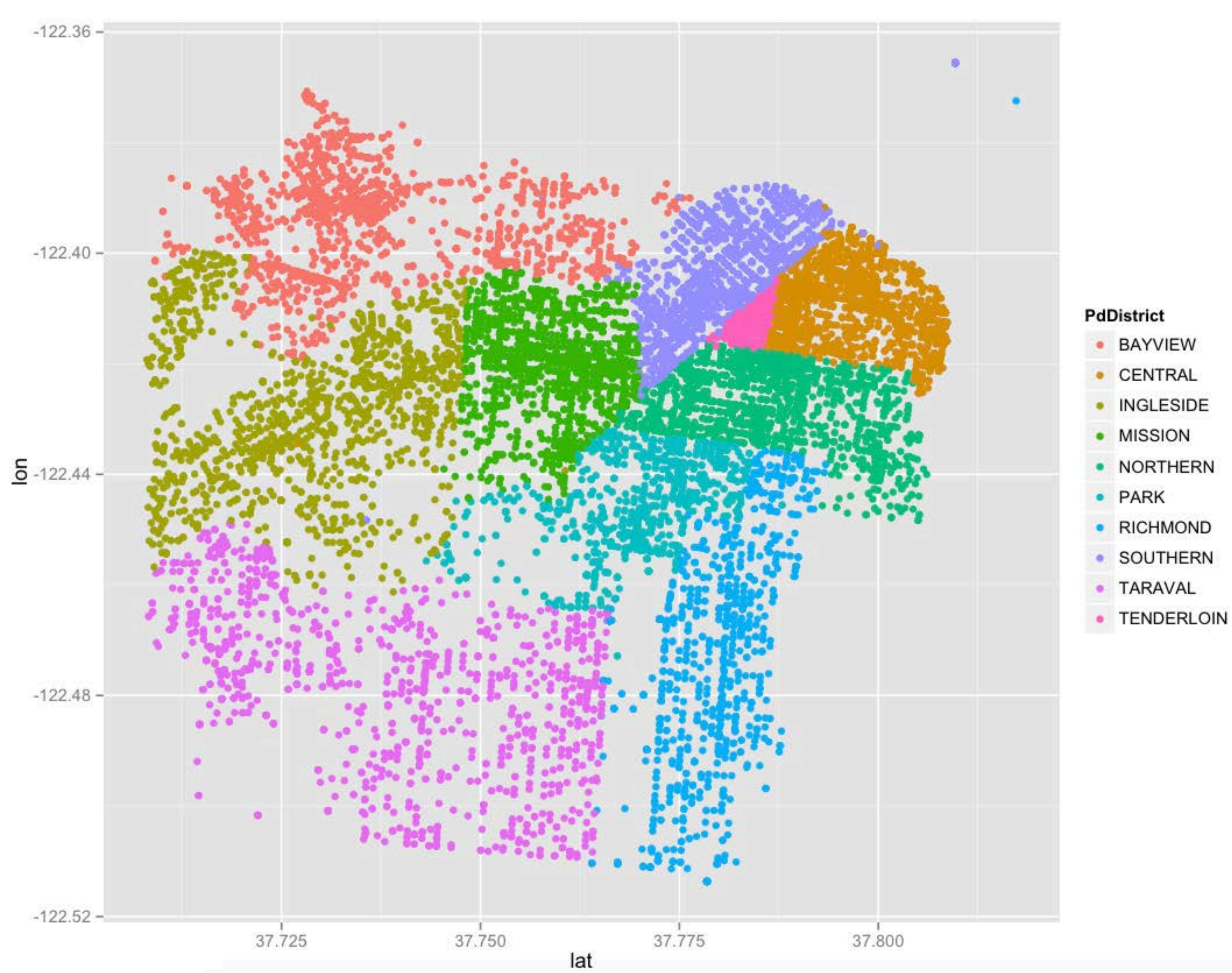
sanfr = read.csv('sanfrancisco_incidents_summer_2014.csv')

summary(sanfr)

library(ggplot2)

#for this I will need latitude/longitude data, luckily it can be easily transformed
#from $Location column
sanfr$lat = as.double(substr(sanfr$Location, 2, 17))
sanfr$lon = as.double(substr(sanfr$Location, 20, 36))

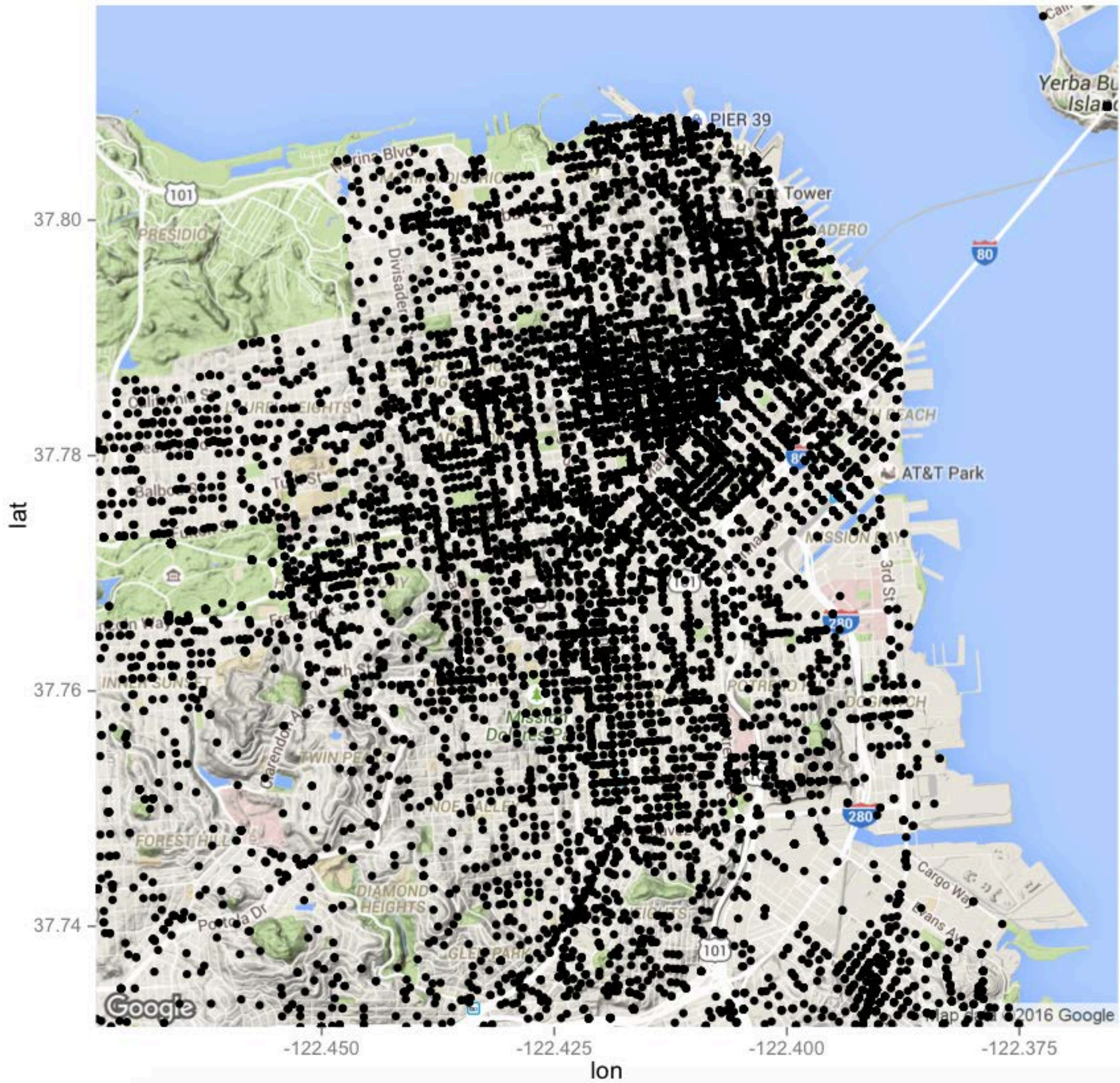
#ok, this works, but maybe seeing it on an actual map would be more intuitive?
p <- ggplot(sanfr, aes(lat, lon, color = PdDistrict))
p + geom_point()
```

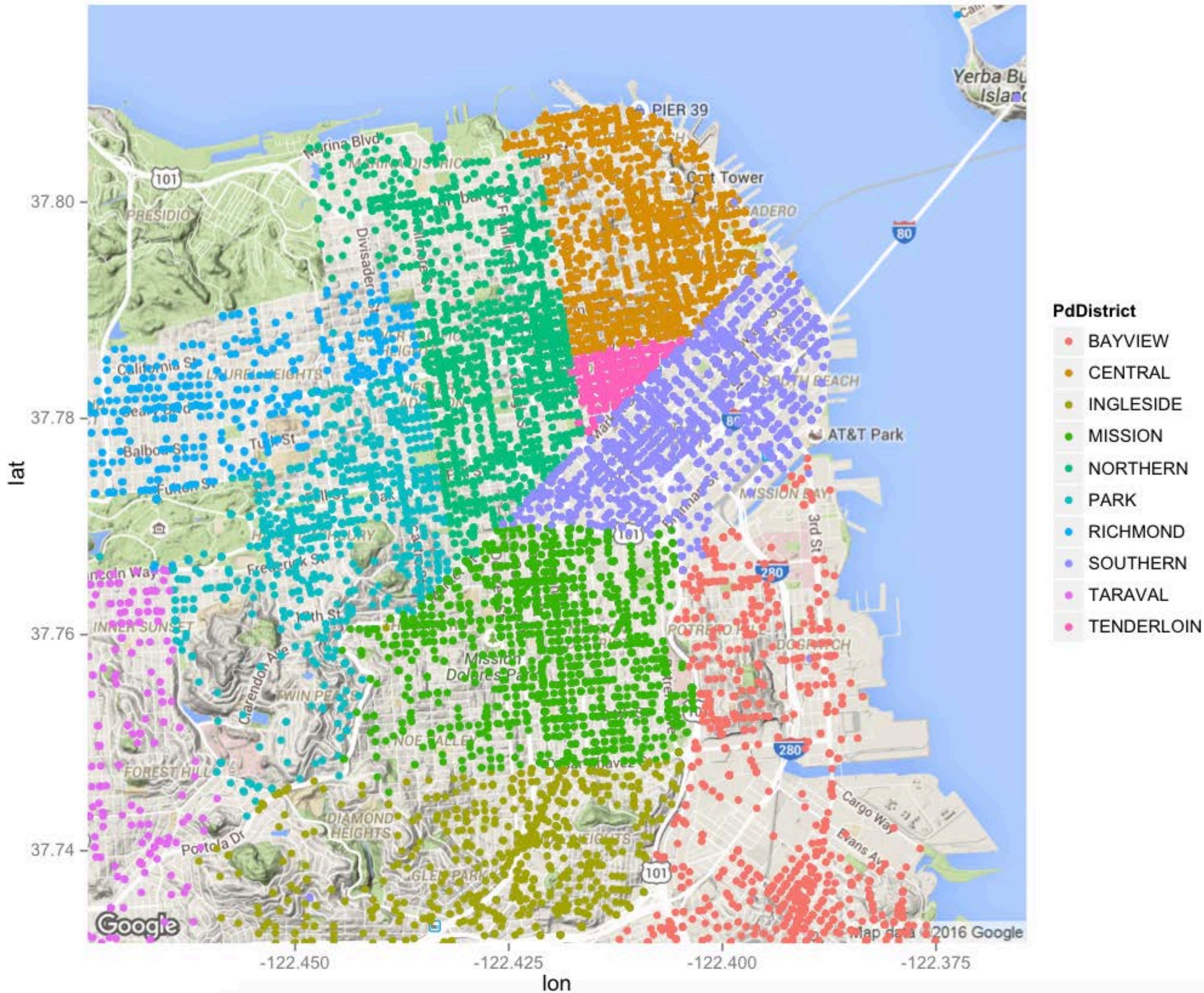
```
library(maps)  
library(ggmap)
```

```
sf = get_map(location = "san francisco", zoom = 13)
```

```
#now it is much more clear where crime happens most  
ggmap(sf) + geom_point(data = sanfr, aes(x = lon, y = lat))
```


```
#and with districts colored it's easier to send police where it's needed  
ggmap(sf) + geom_point(data = sanfr, aes(x = lon, y = lat, color = PdDistrict))
```

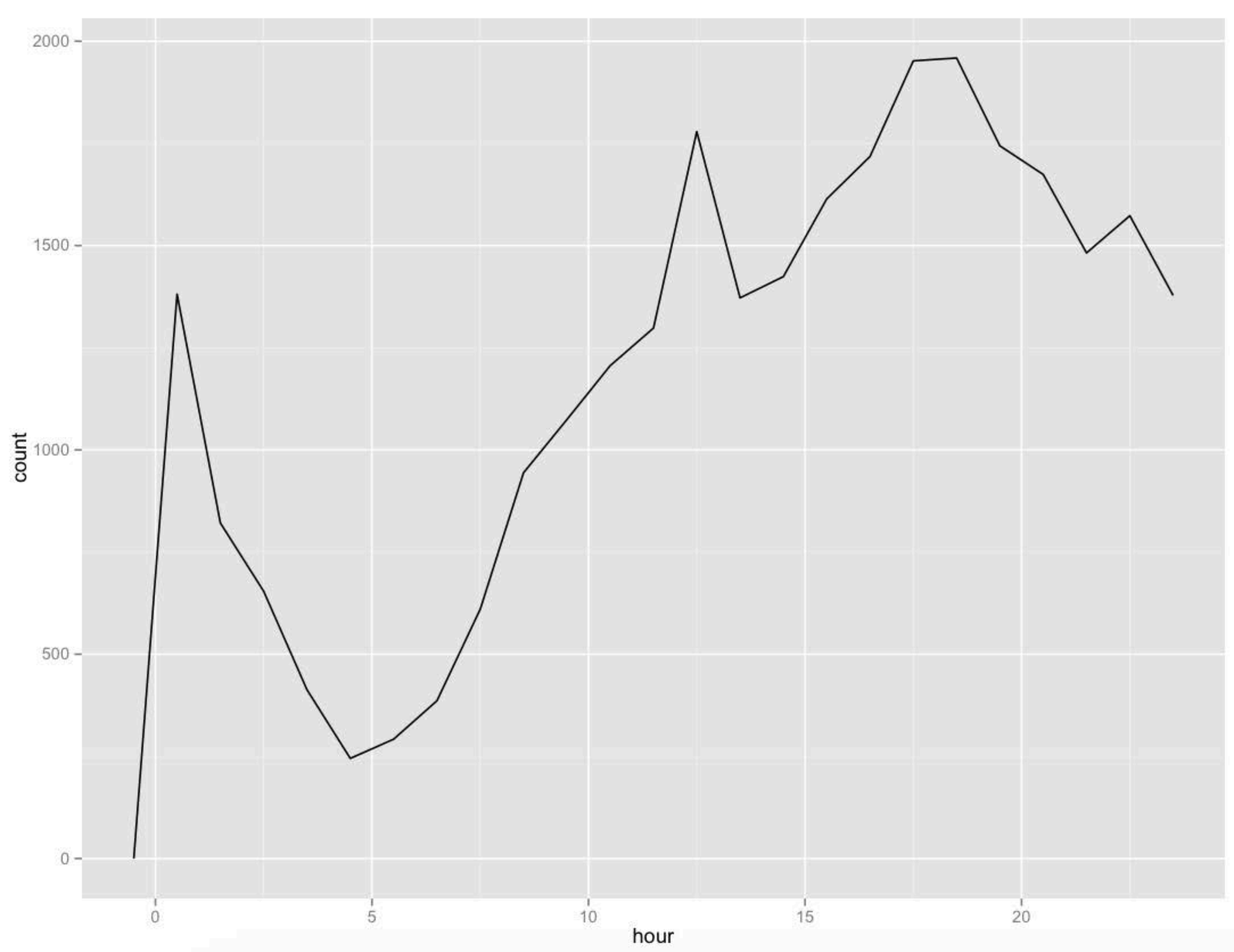

```
#extract hours from $Time column
```

```
sanfr$hour = as.numeric(substr(sanfr$Time, 1, 2))
```

```
#plot it as a line/hist/not sure what better term to use
```

```
#apparently 4AM is the safest hour, did not expect that
```

```
ggplot(sanfr, aes(hour)) + geom_freqpoly(binwidth = 1)
```




```
#histogram of time vs. $DayOfWeek  
#the objective of this one is to identify any outliers  
#Saturday and Sunday at 1AM show about twice as many crimes as other days  
#same thing for midnight on Fridays and Saturdays  
ggplot(sanfr, aes(hour, fill = DayOfWeek)) + geom_histogram(binwidth = 1)
```

