

Crashes in Delaware County in 2016

Graphics with ggplot2

Name: Rajnish Yadav

```
#Put all necessary libraries here
library(ggplot2)
library(readr)
library(dplyr)
library(forcats)
library(viridis)
library(ggthemes)

#Load the dataset
crash_del <- read_csv("~/DataScience/Statistics/Data/crash_del.csv")
#Load uninsured dataset
Uninsured_PA <- read_csv("~/DataScience/Statistics/Data/Uninsured_PA.csv")

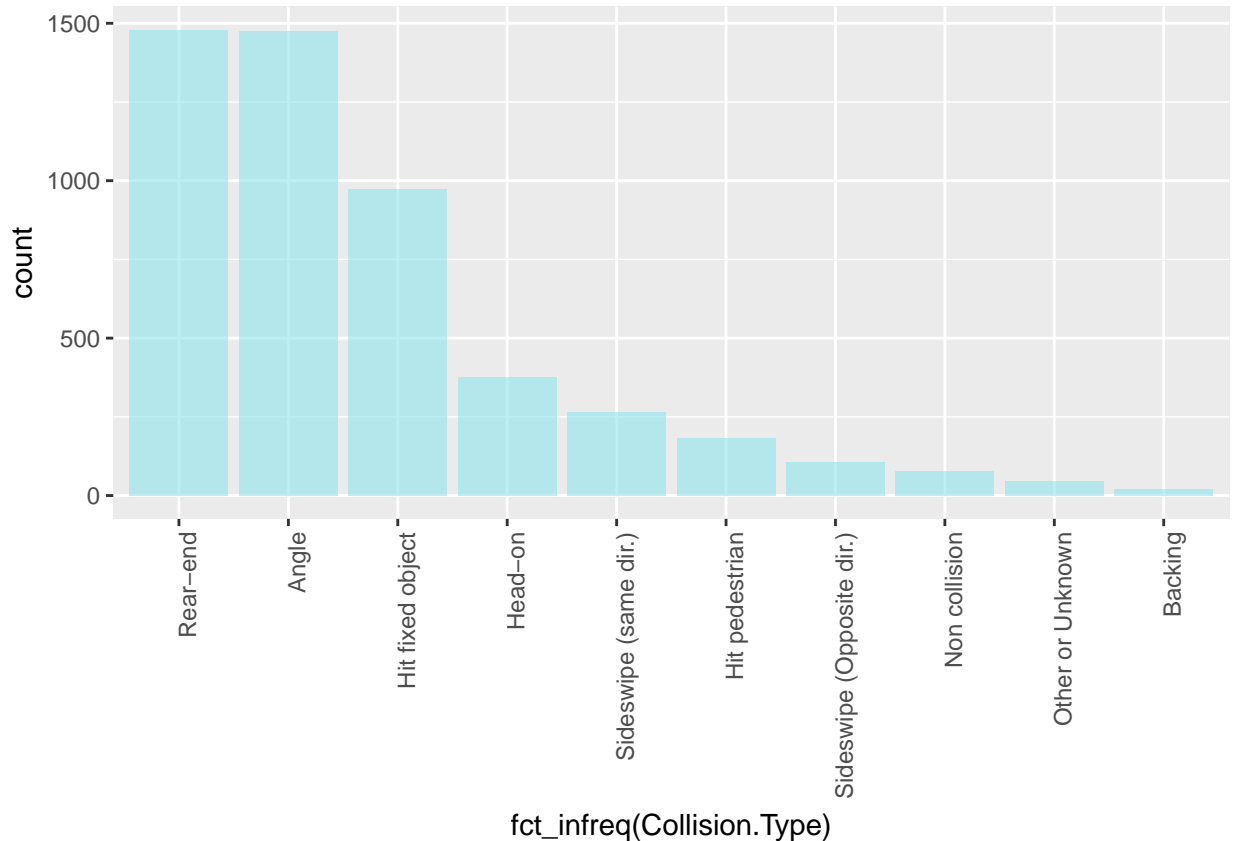
#Fix column names
names(Uninsured_PA) <- gsub(" ", "", names(Uninsured_PA))
names(Uninsured_PA) <- gsub("-", "", names(Uninsured_PA))
```

We will use data on all reported crashes in Delaware County for 2016. These data are provided by the Pennsylvania Department of Transportation. Information the various variables can be found at this [link](#). We will also use data on uninsured rates in PA. Information regarding the various variables can be found at this [link](#).

First of all, let's try to understand the frequency of the different types of collisions (`Collision.Type`).

For ease of reading, we would like to order the bars from most common to least common. In the `forcats` package, there is a function called `fct_infreq`, which will sort the categorical variable by frequency.

```
ggplot(data = crash_del, aes(x = fct_infreq(Collision.Type))) +
  geom_bar(fill = "cadetblue2", alpha = 0.6) +
  theme(axis.text.x=element_text(angle=90, hjust=1))
```



Two most common collision types are 'Rear-end' and 'Angle'. 'Side swipe', 'Non collision', and 'Backing' are the least common collision types. From the bar plot, we can see that 'Rear-end', 'Angle' and 'Hit fixed object' account for ~80% of all collisions.

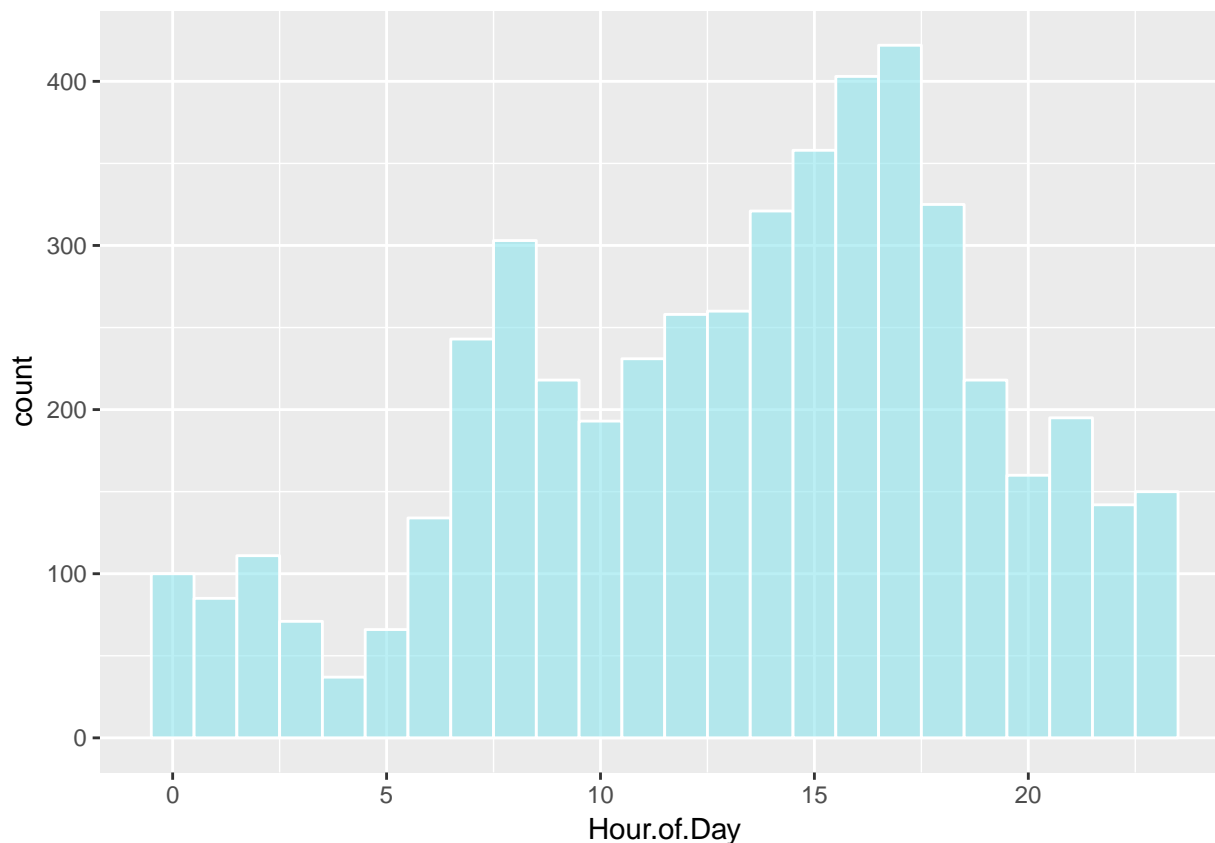
Now we want to plot the distribution of accidents over the hours of the day.

```
crash_del_2 <- filter(crash_del, Hour.of.Day != "99")
```

This code creates a new data frame where any observations where 'Hour.of.Day' = 99 is removed/filtered and this new dataset is assigned to crash_del_2. The size of the new dataset decreased by 3 to 5004 because there were exactly 3 values of 'Hour.of.Day' which were equal to 99.

Visually, some people prefer to see the separation between bars. This can be done by setting the color to white.

```
crash_del_2$Hour.of.Day <- as.integer(crash_del_2$Hour.of.Day)
ggplot(crash_del_2, aes(x = Hour.of.Day)) +
  geom_histogram(binwidth = 1, color = "white", fill = "cadetblue2", alpha = 0.6)
```

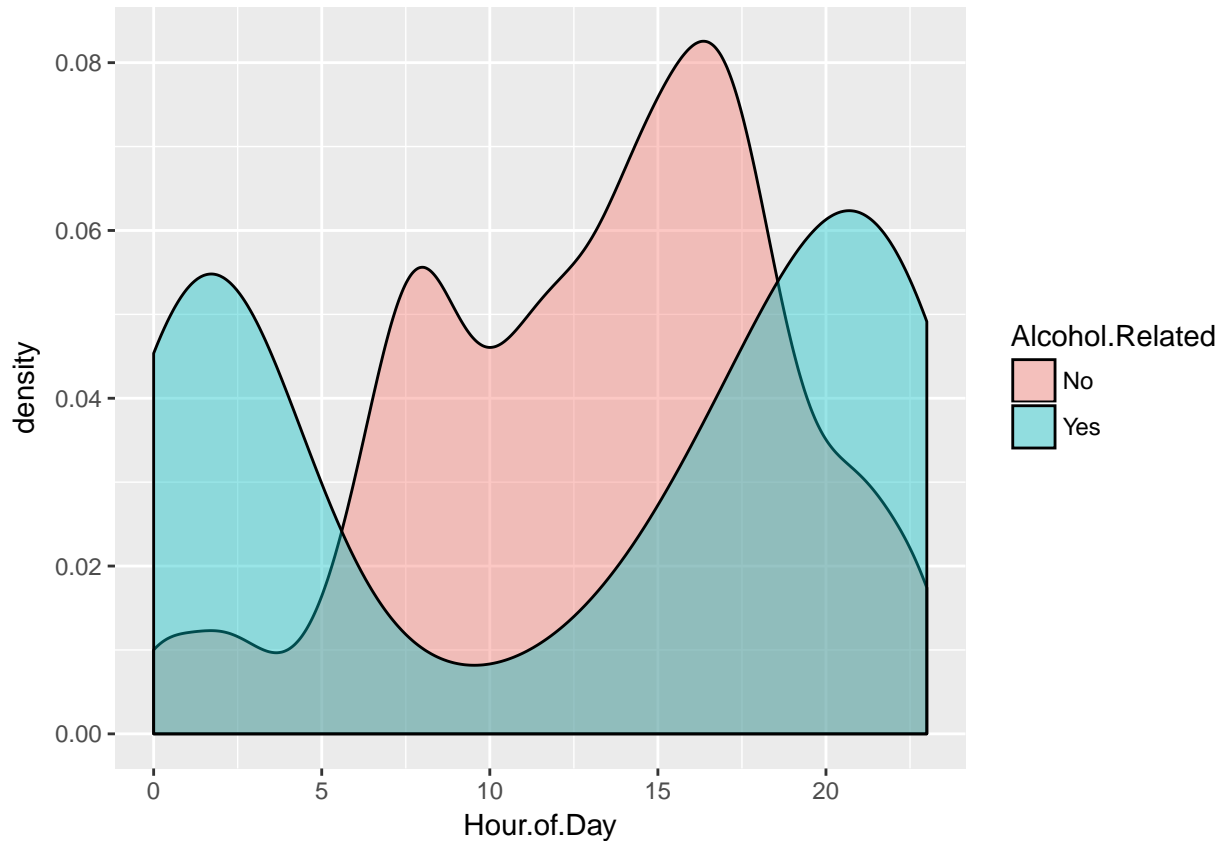


More accidents tend to happen during the morning rush hour (7 am to 8 am) and the afternoon rush hour (4 pm to 5 pm). However, the afternoon rush hour has more number of accidents compared to the morning rush hour. Not many people after midnight, so naturally less number of accidents occur from midnight to 5 am.

Note that we cannot conclude that it is safest to drive late at night because we do not know about the proportions of accidents occurring at night i.e there might be less cars on the road at night resulting into less number of accidents at any given hour. Moreover, it is also possible that more people might be driving under the influence of alcohol resulting into more number of accidents per 100 cars on the road. Thus, it is not a good idea to assume that it is safest to drive late at night.

Now, let's explore whether the distribution of accidents over the course of a day vary if the accident is alcohol related. We will use `crash_del_2`.

```
ggplot(crash_del_2, aes(x = Hour.of.Day, fill = Alcohol.Related)) +  
  geom_density(alpha = 0.4)
```



Accidents related to alcohol occur more from midnight to about 2 am at night. One reason might be that people are driving back home from bars and parties at night under the influence. The number of alcohol related accidents start increasing from afternoon approx. 5 pm until night approx. 9 pm. Less alcohol related accidents occur during the day. More non-alcohol accidents occur during the day.

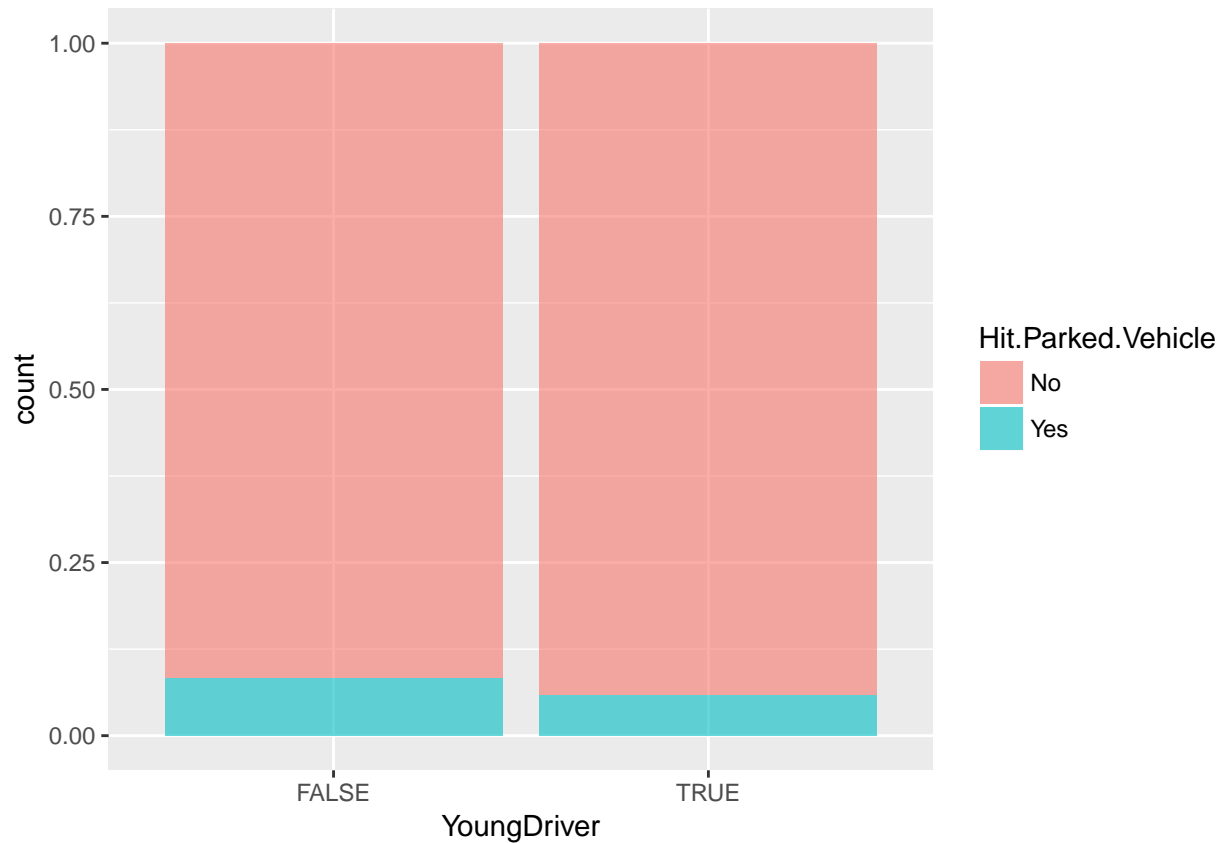
Young drivers and old drivers get a bad rap. Let's explore some accident trends associated with these types of drivers.

```
crash_del <- mutate(crash_del, YoungDriver = (Driver.16.Years.Old == "Yes" | Driver.17.Years.Old=="Yes"))
```

The above code is adding two extra columns "YoungDriver" and "OlderDriver" to the crash_del dataset. 'YoungDriver' refer to the younger drivers who are between ages 16 and 18 inclusive. 'OlderDriver' refers to the older drivers who are older than or equal to 65 years. This adds the number of variables by 2 but the number of observations remain unchanged.

The plot below answers the following question: Given an accident occurred, are younger drivers more likely to hit a parked vehicle than non-young drivers?

```
ggplot(data = crash_del, aes(x = YoungDriver, fill = Hit.Parked.Vehicle)) +  
  geom_bar(position = "fill", alpha = 0.6)
```



Based on the bar plot above, we can say that given an accident occurred, younger drivers between the ages of 16 and 18 inclusive are not more likely to hit a parked vehicle than non-young drivers.

The plot below answers the question posed above but this time for older drivers.

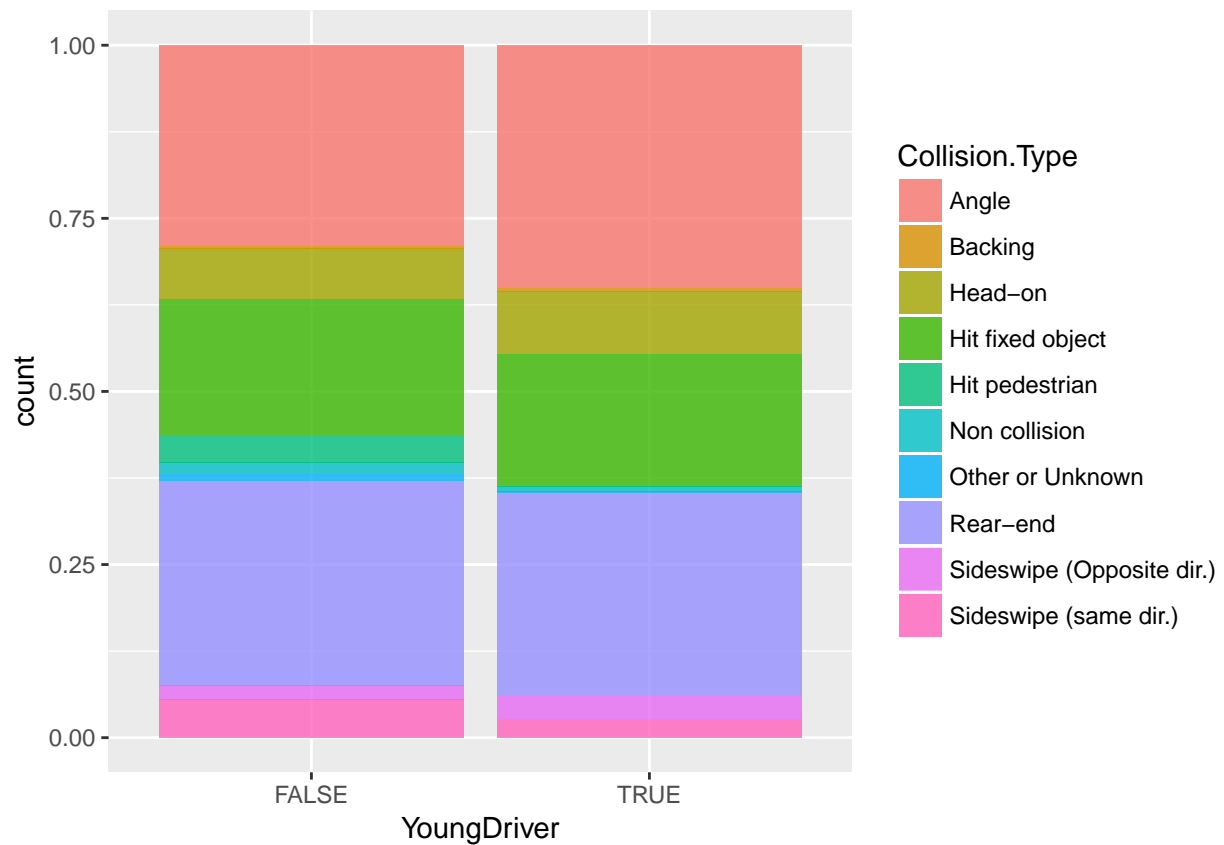
```
ggplot(data = crash_del, aes(x = OlderDriver, fill = Hit.Parked.Vehicle)) +  
  geom_bar(position = "fill", alpha = 0.6)
```



Based on the bar plot above, we can say that given an accident occurred, older drivers above age 65 (inclusive) are not more likely to hit a parked vehicle than drivers below age 65.

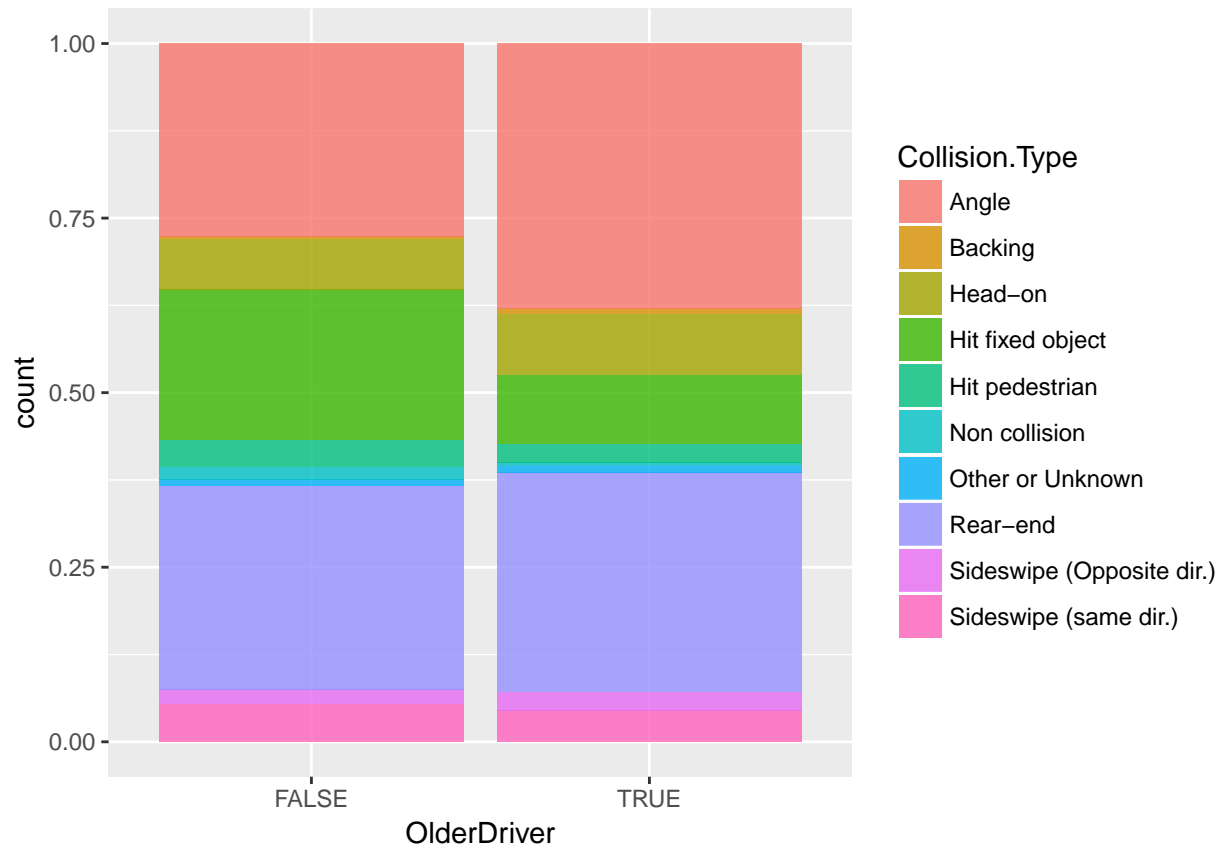
Let's create a barplot that displays the conditional proportions of the collision types for young and non-young drivers.

```
ggplot(crash_del, aes(x = YoungDriver, fill = Collision.Type)) +  
  geom_bar(position = "fill", alpha = 0.8)
```



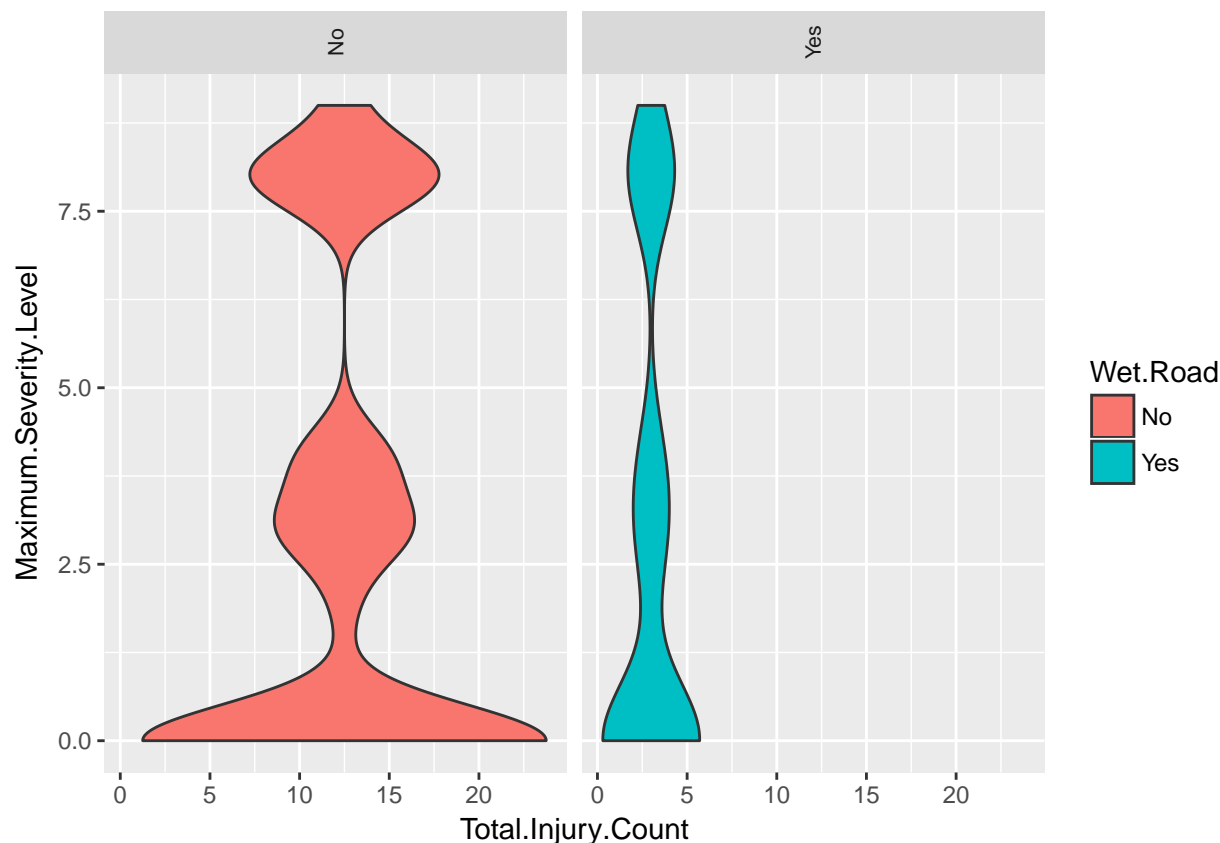
Let's create a barplot that displays the conditional proportions of the collision types for older and non-older drivers.

```
ggplot(crash_del, aes(x = OlderDriver, fill = Collision.Type)) +
  geom_bar(position = "fill", alpha = 0.8)
```



Based on your plots, 'Angle' collisions, 'Head-on' and 'Side swipe(Opposite dir.)' collisions are more common when a younger driver is involved compared to non-young drivers. 'Angle' collisions and 'Rear-end' collisions are more common when older drivers are involved compared to non-older drivers.

```
ggplot(crash_del, aes(x = Total.Injury.Count, y = Maximum.Severity.Level, fill = Wet.Road)) +
  geom_violin() +
  facet_grid(. ~ Wet.Road) +
  theme(strip.text.x = element_text(size = 8, angle = 90))
```

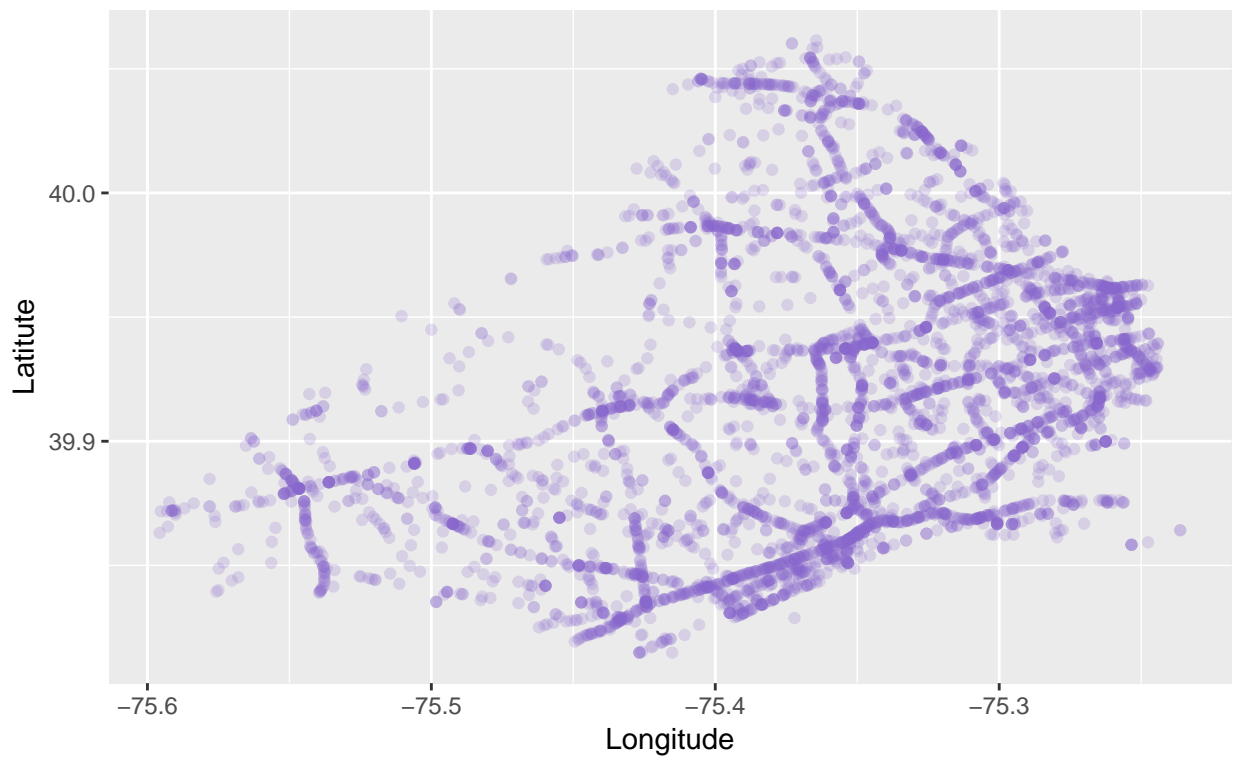
Three variables are Total.Injury.Count (number of people getting injured in accidents), Maximum.Severity.Level (injury severity level; larger number implies more severe injury) and Wet.Road (whether the roads are wet or not). I would expect people to drive slowly when the road is wet. Hence, if an accident occurs less number of people are injured and their injury severity level is from small to severe; however, most people are not severely injured. When the road is not wet, people drive faster and might speed at times, resulting into more number of people getting more severely injured.

Let's explore the locations of crashes in Delaware County.

```
ggplot(crash_del, aes(x = Longitude..Decimal., y = Latitude..Decimal.)) +
  geom_point(alpha = 0.2, color = "mediumpurple3") +
  labs(title = "Scatterplot of locations of crashes in Delaware County",
       caption = "Data Source: Pennsylvania Department of Transportation",
       x = "Longitude", y = "Latitude")
```

```
## Warning: Removed 50 rows containing missing values (geom_point).
```

Scatterplot of locations of crashes in Delaware County



Data Source: Pennsylvania Department of Transportation

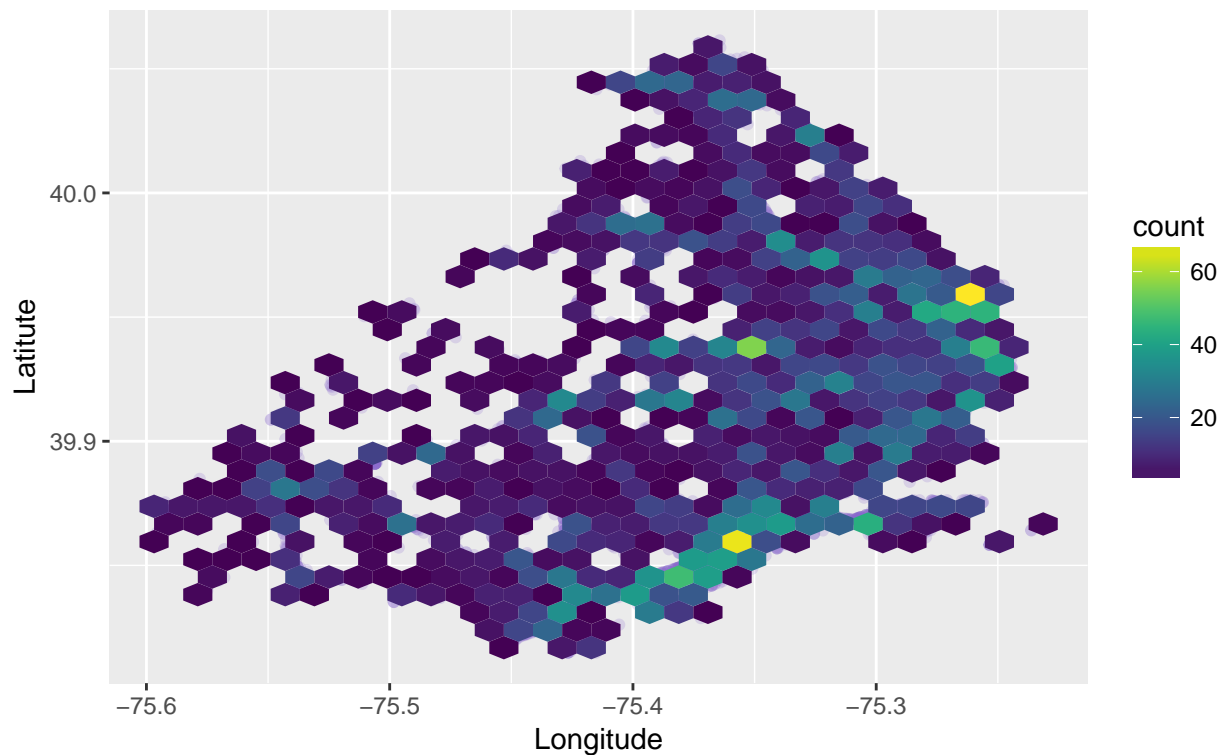
Let's create a heatmap.

```
ggplot(crash_del, aes(x = Longitude..Decimal., y = Latitude..Decimal.)) +  
  geom_point(alpha = 0.2, color = "mediumpurple3") +  
  labs(title = "Heatmap of locations of crashes in Delaware County",  
       caption = "Data Source: Pennsylvania Department of Transportation",  
       x = "Longitude", y = "Latitude") +  
  geom_hex() +  
  scale_fill_viridis()
```

```
## Warning: Removed 50 rows containing non-finite values (stat_binhex).
```

```
## Warning: Removed 50 rows containing missing values (geom_point).
```

Heatmap of locations of crashes in Delaware County



Data Source: Pennsylvania Department of Transportation

Let's re-create the scatterplot and heatmap but this time point out where Swarthmore, PA is located.

```
swarthmore <- filter(crash_del, Municipality.Name=="Swarthmore Boro") %>%
  select(Longitude..Decimal., Latitude..Decimal.)
swarthmore
```

```
## # A tibble: 43 x 2
##   Longitude..Decimal. Latitude..Decimal.
##   <dbl>             <dbl>
## 1      -75.3         39.9
## 2      -75.3         39.9
## 3      -75.4         39.9
## 4      -75.4         39.9
## 5      -75.3         39.9
## 6      -75.4         39.9
## 7      -75.3         39.9
## 8      -75.3         39.9
## 9      -75.3         39.9
## 10     -75.4         39.9
## # ... with 33 more rows
```

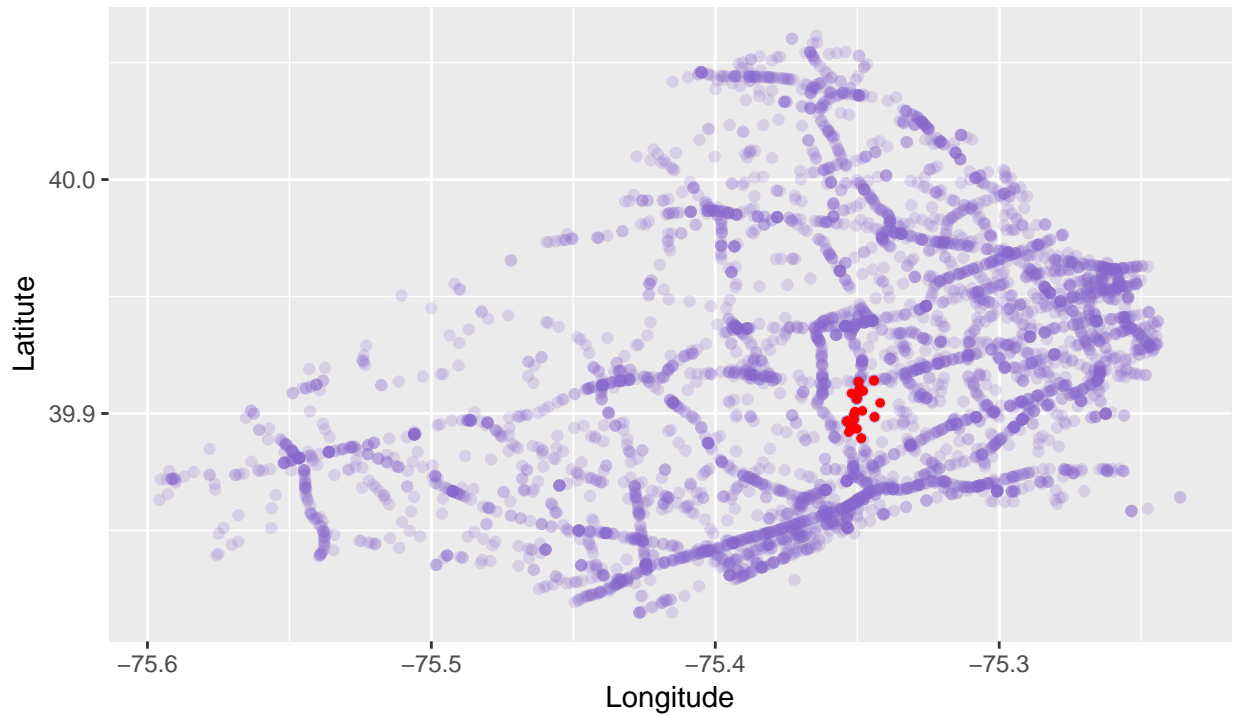
```
ggplot(crash_del, aes(x = Longitude..Decimal., y = Latitude..Decimal.)) +
  geom_point(alpha = 0.2, color = "mediumpurple3") +
  geom_point(data = swarthmore, aes(x = Longitude..Decimal., y = Latitude..Decimal.), color = "red", size = 100) +
  labs(title = "Scatterplot of locations of crashes in Delaware County",
        subtitle = "Red represents where Swarthmore is located",
        caption = "Data Source: Pennsylvania Department of Transportation",
```

```
x = "Longitude", y = "Latitude")
```

```
## Warning: Removed 50 rows containing missing values (geom_point).
```

Scatterplot of locations of crashes in Delaware County

Red represents where Swarthmore is located

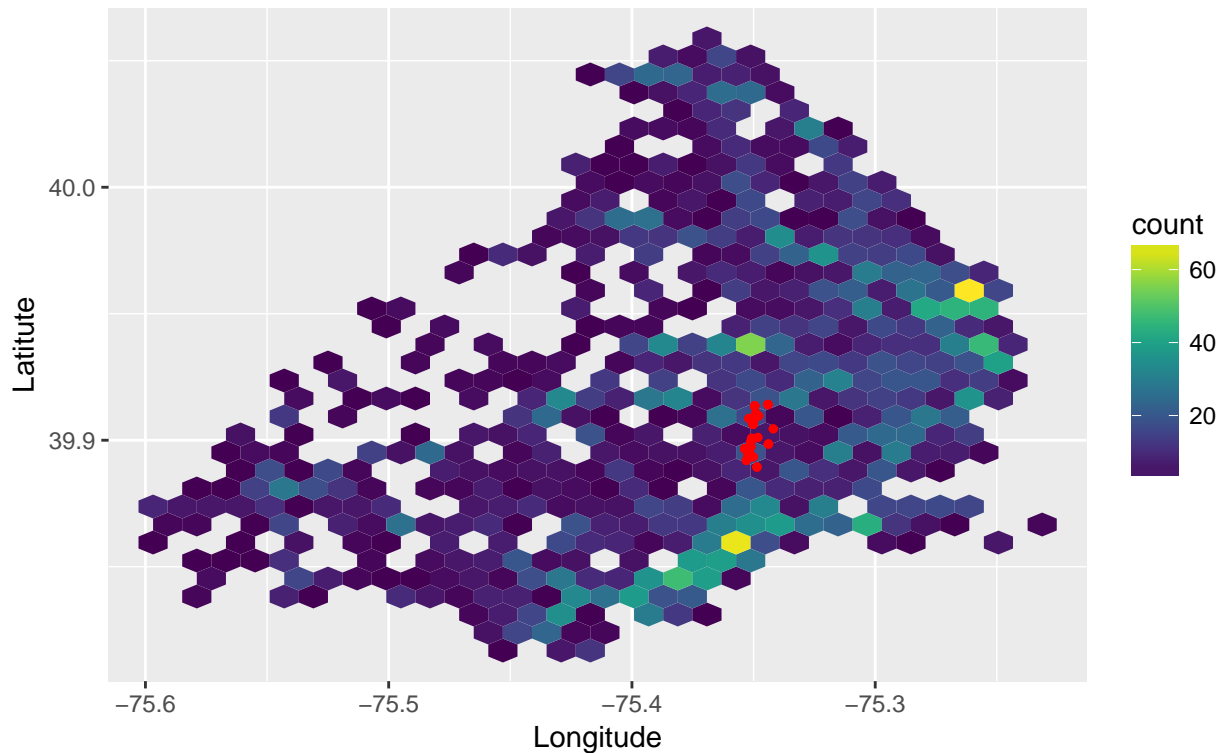


Data Source: Pennsylvania Department of Transportation

```
ggplot(crash_del, aes(x = Longitude..Decimal., y = Latitude..Decimal.)) +
  geom_hex() +
  geom_point(data = swarthmore, aes(x = Longitude..Decimal., y = Latitude..Decimal.), color = "red", size = 100) +
  labs(title = "Heatmap of locations of crashes in Delaware County",
       caption = "Data Source: Pennsylvania Department of Transportation; red represents where Swarthmore is located",
       x = "Longitude", y = "Latitude") +
  scale_fill_viridis()
```

```
## Warning: Removed 50 rows containing non-finite values (stat_binhex).
```

Heatmap of locations of crashes in Delaware County

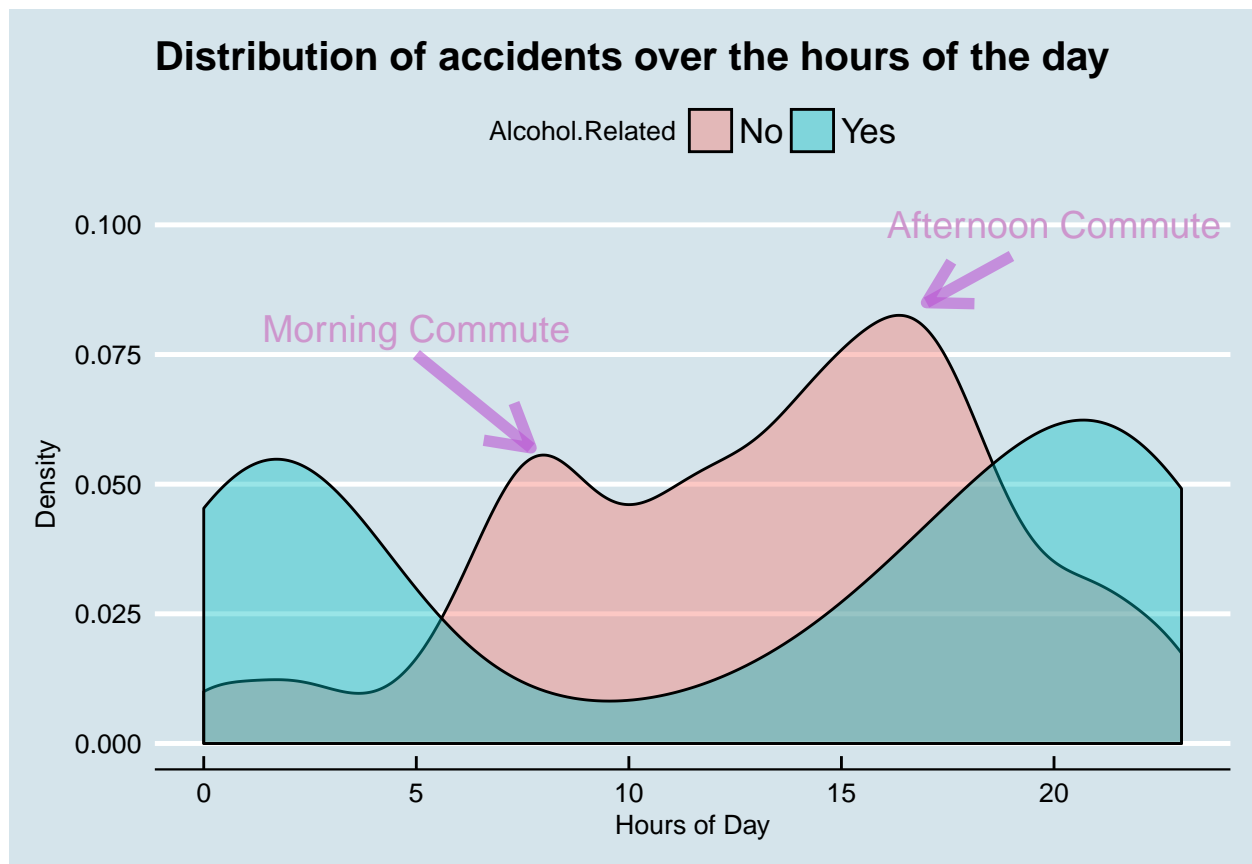


œ: Pennsylvania Department of Transportation; red represents where Swarthmore is located

The crashes seem to be clumping around the I-95 and I-476 highways which pass through the Delaware County. The scatterplot with heatmap gives us a better sense of where the crashes are occurring.

Let's customize and add context to the graphic we created earlier. Recall that we created a faceted density plot of the distribution of accidents over the hours of the day by whether or not it was alcohol related.

```
crash_del$Hour.of.Day <- as.integer(crash_del$Hour.of.Day)
crash_del <- filter(crash_del, Hour.of.Day != "99")
ggplot(crash_del, aes(x = Hour.of.Day, fill = Alcohol.Related)) +
  geom_density(alpha = 0.4) +
  labs(title = "Distribution of accidents over the hours of the day", x = "Hours of Day", y = "Density") +
  annotate("segment", x = 5, xend = 7.7, y = 0.075,
    yend = 0.057, colour = "mediumorchid", size = 2, alpha = .6,
    arrow = arrow()) + annotate("segment", x = 19,
    xend = 17, y = 0.094, yend = 0.085, colour = "mediumorchid",
    size = 2, alpha = 0.6, arrow = arrow()) + annotate("text",
    x = 5, y = 0.08, label = c("Morning Commute"),
    color = "plum3", size = 5) + annotate("text", x = 20,
    y = 0.1, label = c("Afternoon Commute"), color = "plum3",
    size = 5) +
  theme_economist()
```



Let's look at the relationship between household income and the percentage of the county that is uninsured.

```
Uninsured_PA <- mutate(Uninsured_PA, UninsuredPercent = Uninsured/Population)
```

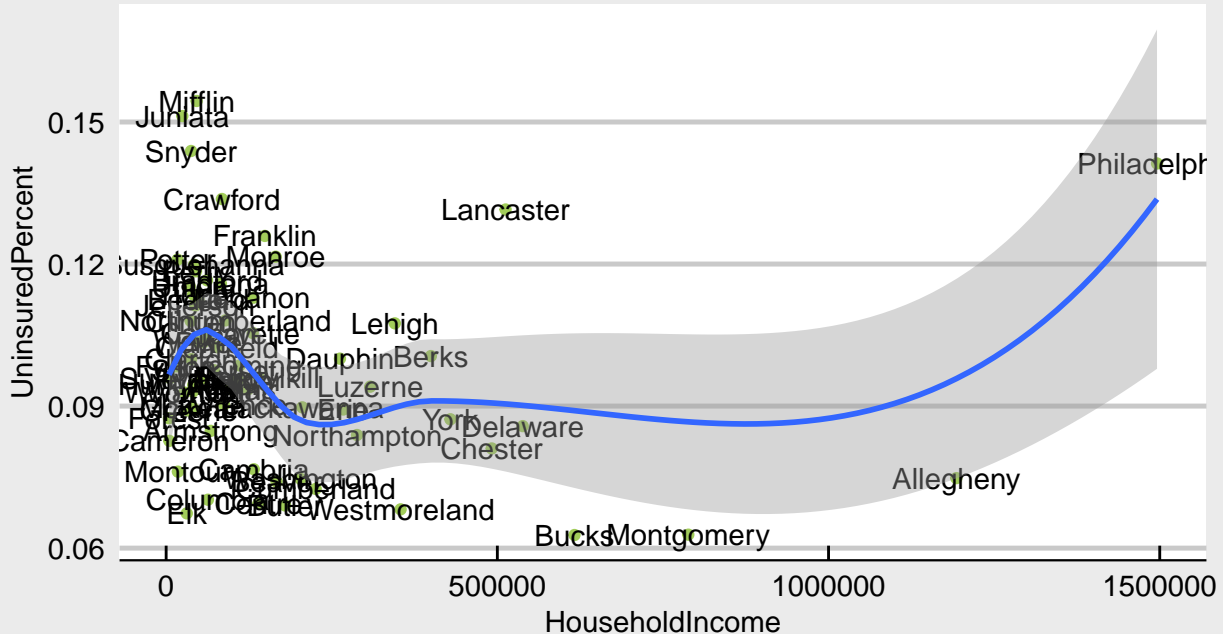
We are computing a new column called UninsuredPercent by dividing the Uninsured with the Population. We are then appending this new column to Uninsured_PA.

Let's create a scatterplot. Include a smoother and the county names.

```
ggplot(Uninsured_PA, aes(x = HouseholdIncome, y = UninsuredPercent)) +
  geom_point(color = "darkolivegreen3") +
  geom_text(aes(label = County)) +
  geom_smooth() +
  labs(title = "Relationship between household income and the percentage
of the county that is uninsured", subtitle = "Data Provided by
Small Area Health Insurance Estimates Program, U.S. Census Bureau") +
  theme_economist_white()
```

Relationship between household income and the percentage of the county that is uninsured

Data Provided by
Small Area Health Insurance Estimates Program, U.S. Census Bureau



The relationship between the variables is difficult to tell from the graph above. Many counties are concentrated in the low income household distribution with 9-12 percentage of uninsured population. The curve is skewed and it looks like, in general, population with higher household income has lower uninsured percentage.

Philadelphia and Allegheny seem to be the outlier. We see extremely high income level for Philadelphia and Allegheny; moreover, Philadelphia has a high uninsured rate. I think the sample might not have been taken randomly for Philadelphia or there is some mistake in entering the data for household income for Philadelphia in the dataset. Lancaster seems to be a rich county based on the household income, however it does have a high insured rate (>12%).

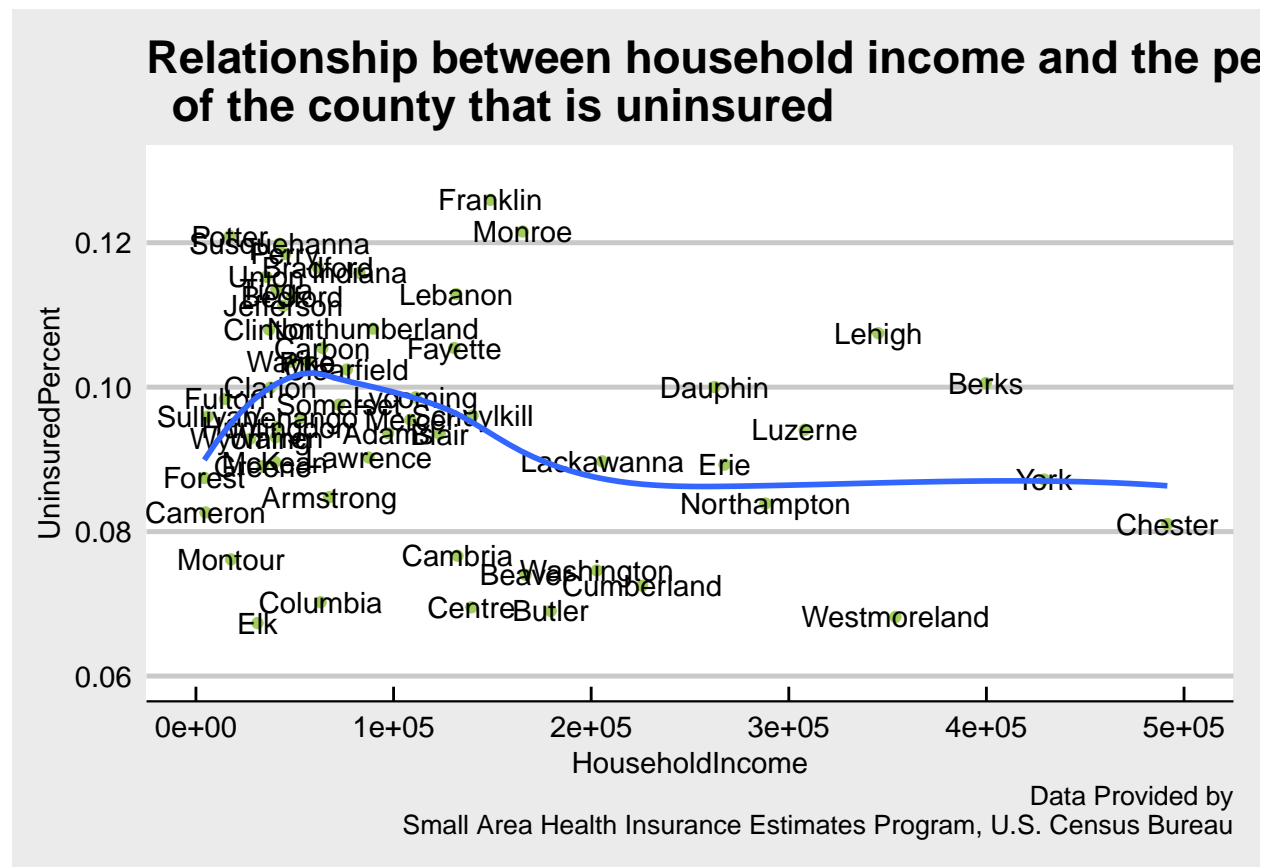
Recreate the scatterplot but this time modify the scale of the axes so that the outliers are not outlying.

```
ggplot(Uninsured_PA, aes(x = HouseholdIncome, y = UninsuredPercent)) +
  geom_point(color = "darkolivegreen3") +
  geom_text(aes(label = County)) +
  geom_smooth(se = FALSE) +
  xlim(0, 500000) +
  ylim(0.06, 0.13) +
  labs(title = "Relationship between household income and the percentage
of the county that is uninsured", caption = "Data Provided by
Small Area Health Insurance Estimates Program, U.S. Census Bureau") +
  theme_economist_white()
```

Warning: Removed 10 rows containing non-finite values (stat_smooth).

```
## Warning: Removed 10 rows containing missing values (geom_point).
```

```
## Warning: Removed 10 rows containing missing values (geom_text).
```



We can also transform the coordinate system instead of the scale using the following layer. Let's modify the graph so that we are transforming the coordinate system instead of the scale.

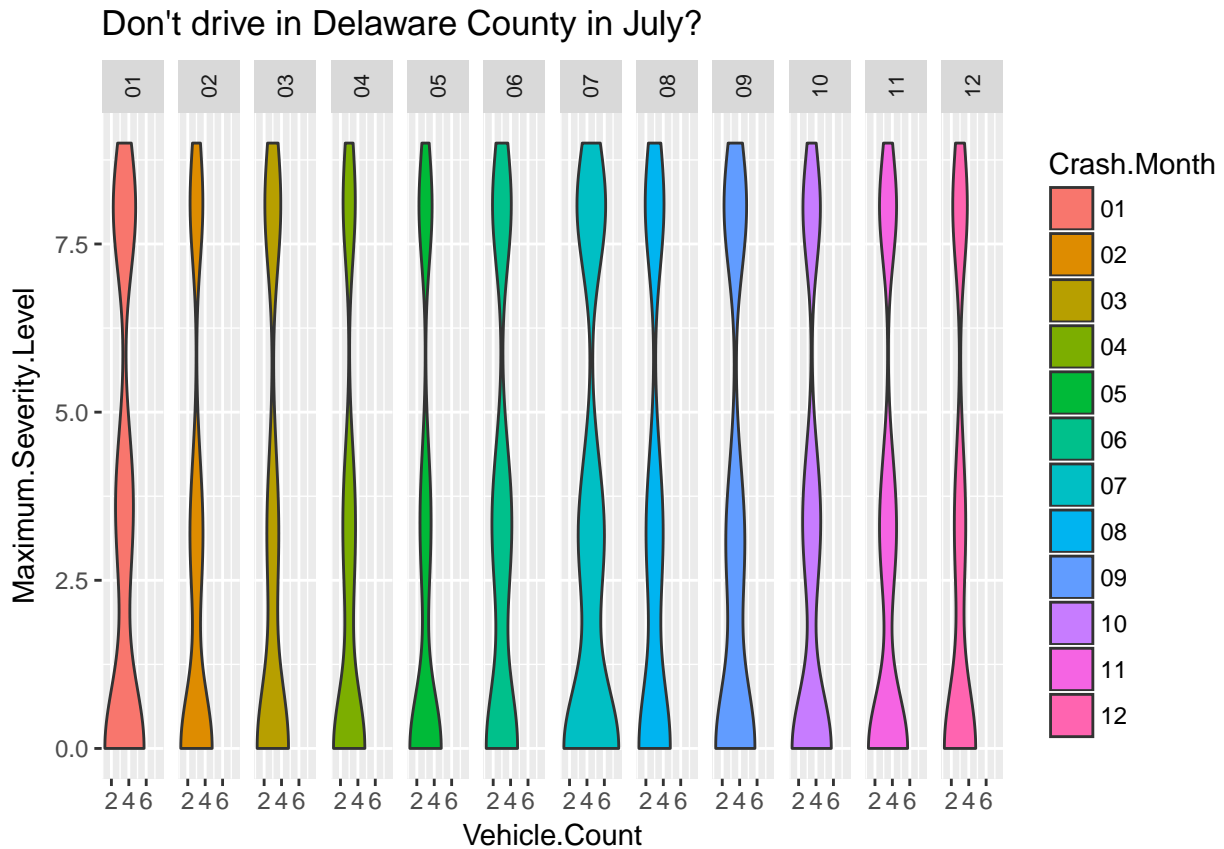
```
ggplot(Uninsured_PA, aes(x = HouseholdIncome, y = UninsuredPercent)) +  
  geom_point(color = "darkolivegreen3") +  
  geom_text(aes(label = County)) +  
  geom_smooth(se = FALSE) +  
  labs(title = "Relationship between household income and the percentage  
of the county that is uninsured", caption = "Data Provided by  
Small Area Health Insurance Estimates Program, U.S. Census Bureau") +  
  theme_economist_white() +  
  coord_trans(x = "log10")
```


This scatter plot displays the relationship between household income and the percentage of the uninsured population across Pennsylvania's counties. The x-axis, 'HouseholdIncome', is on a logarithmic scale ranging from approximately 100,000 to over 1,000,000. The y-axis, 'UninsuredPercent', ranges from 0.08 to 0.14. A blue trend line shows a general downward trend, indicating that as household income increases, the percentage of the uninsured population tends to decrease. The data points are labeled with county names, showing a wide distribution of uninsured rates across different income levels.

County	Household Income (approx.)	Uninsured Percent (approx.)
Forest	100,000	0.085
Cameron	100,000	0.082
Sullivan	100,000	0.095
Fulton	200,000	0.098
Montour	200,000	0.075
Potter	250,000	0.120
Wyo	250,000	0.090
Greene	300,000	0.088
Marion	300,000	0.085
Armstrong	350,000	0.082
Columbia	400,000	0.075
Elk	300,000	0.070
Lawrence	400,000	0.088
Adams	450,000	0.090
Franklin	500,000	0.125
Monroe	550,000	0.120
Lebanon	550,000	0.110
Fayette	550,000	0.105
Dauphin	600,000	0.095
Berks	650,000	0.098
Luzerne	650,000	0.090
Northampton	700,000	0.085
York	750,000	0.088
Ches	800,000	0.082
Allegheny	1,000,000	0.075
Block	800,000	0.070
Moreland	850,000	0.072
Phila	1,200,000	0.140

The above graph has x-axis with a logarithmic scale, so the values we see are the $\log_{10}()$ values of the household income whereas in the graph before that, the household income refers to the actual household income. Taking $\log_{10}()$ of the household income distribution distributes the x-coordinates more evenly compared to the one without logarithmic x-axis.

17



From the above plot, we can clearly see that July stands out as the month in which the highest number of vehicle accidents occur and people get severely injured in those accidents. Similarly, December and May seem to have the lowest number of vehicle accidents. In each month, there are vehicle accidents ranging from the least severe to the most severe.