# NYPD Shooting Incident Data Report

2024-07-16

## Synopsis

This report analyzes a dataset from the City of New York, retrieved from following data source URL, https: //data.cityofnewyork.us/api/views/833y-fsy8/rows.csv?accessType=DOWNLOAD. The analysis includes data loading, cleaning, exploratory data analysis (EDA), and visualization to uncover insights related to NYPD Shooting Incident Data (Historic).It includes series of steps like Import, tidy and analyze the NYPD Shooting Incident dataset obtained. Also ensure project is reproducible and contains some visualization and analysis. Also includes bias in the data and conclusion of my analysis.

## Step 1: Import Library

```
# Install required packages if not already installed
# install.packages("tidyverse")
# install.packages("lubridate")

library(tidyverse)
library(lubridate)
```

## Step 2: Load Data

# Load the dataset from the provided URL

```
df_nypd = read_csv("https://data.cityofnewyork.us/api/views/833y-fsy8/rows.csv?accessType=DOWNLOAD")
```

```
## Rows: 28562 Columns: 21
## -- Column specification ---------------------------------------------------------
## Delimiter: ","
## chr  (12): OCCUR_DATE, BORO, LOC_OF_OCCUR_DESC, LOC_CLASSFCTN_DESC, LOCATION...
## dbl   (7): INCIDENT_KEY, PRECINCT, JURISDICTION_CODE, X_COORD_CD, Y_COORD_CD...
## lgl   (1): STATISTICAL_MURDER_FLAG
## time  (1): OCCUR_TIME
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
head(df)
```

```
## 
## 1 function (x, df1, df2, ncp, log = FALSE)
## 2 {
## 3     if (missing(ncp))
## 4         .Call(C_df, x, df1, df2, log)
## 5     else .Call(C_dnf, x, df1, df2, ncp, log)
## 6 }
```

## Step 3: Tidy and Transform Data

Let's first eliminate the columns I do not need for this analysis, which are: **PRECINCT**,**JURISDICTION_CODE**,**LOCA
**X_COORD_CD**, **Y_COORD_CD**, and **Lon_Lat**.

```
df_nypd_2 = df_nypd %>% select(INCIDENT_KEY,
                    OCCUR_DATE,
                    OCCUR_TIME,
                    BORO,
                    STATISTICAL_MURDER_FLAG,
                    PERP_AGE_GROUP,
                    PERP_SEX,
                    PERP_RACE,
                    VIC_AGE_GROUP,
                    VIC_SEX,
                    VIC_RACE,
                    Latitude,
                    Longitude)

# Return the column name along with the missing values
lapply(df_nypd_2, function(x) sum(is.na(x)))
```

```
## $INCIDENT_KEY
## [1] 0
## 
## $OCCUR_DATE
## [1] 0
## 
## $OCCUR_TIME
## [1] 0
## 
## $BORO
## [1] 0
## 
## $STATISTICAL_MURDER_FLAG
## [1] 0
## 
## $PERP_AGE_GROUP
## [1] 9344
## 
## $PERP_SEX
## [1] 9310
## 
## $PERP_RACE
## [1] 9310
```

```
##
## $VIC_AGE_GROUP
## [1] 0
##
## $VIC_SEX
## [1] 0
##
## $VIC_RACE
## [1] 0
##
## $Latitude
## [1] 59
##
## $Longitude
## [1] 59
```

Understanding the reasons why data are missing is important for handling the remaining data correctly.
There's a fair amount of unidentifiable data on perpetrators (age, race, or sex.) Those cases are possibly still
active and ongoing investigation. In fear of missing meaningful information, I handle this group of missing
data by calling them as another group of "Unknown".

Key observations on data type conversion are:

- **INCIDENT_KEY** should be treated as a string.
- **BORO** should be treated as a factor.
- **PERP_AGE_GROUP** should be treated as a factor.
- **PERP_SEX** should be treated as a factor.
- **PERP_RACE** should be treated as a factor.
- **VIC_AGE_GROUP** should be treated as a factor.
- **VIC_SEX** should be treated as a factor.
- **VIC_RACE** should be treated as a factor.

```r
# Tidy and transform data
df_nypd_2 = df_nypd_2 %>%
  replace_na(list(PERP_AGE_GROUP = "Unknown", PERP_SEX = "Unknown", PERP_RACE = "Unknown"))

# Remove extreme values in data
df_nypd_2 = subset(df_nypd_2, PERP_AGE_GROUP!="1020" & PERP_AGE_GROUP!="224" & PERP_AGE_GROUP!="940")

df_nypd_2$PERP_AGE_GROUP = recode(df_nypd_2$PERP_AGE_GROUP, UNKNOWN = "Unknown")
df_nypd_2$PERP_SEX = recode(df_nypd_2$PERP_SEX, U = "Unknown")
df_nypd_2$PERP_RACE = recode(df_nypd_2$PERP_RACE, UNKNOWN = "Unknown")
df_nypd_2$VIC_SEX   = recode(df_nypd_2$VIC_SEX, U = "Unknown")
df_nypd_2$VIC_RACE    = recode(df_nypd_2$VIC_RACE, UNKNOWN = "Unknown")
df_nypd_2$INCIDENT_KEY = as.character(df_nypd_2$INCIDENT_KEY)
df_nypd_2$BORO = as.factor(df_nypd_2$BORO)
df_nypd_2$PERP_AGE_GROUP = as.factor(df_nypd_2$PERP_AGE_GROUP)
df_nypd_2$PERP_SEX = as.factor(df_nypd_2$PERP_SEX)
df_nypd_2$PERP_RACE = as.factor(df_nypd_2$PERP_RACE)
df_nypd_2$VIC_AGE_GROUP = as.factor(df_nypd_2$VIC_AGE_GROUP)
df_nypd_2$VIC_SEX = as.factor(df_nypd_2$VIC_SEX)
df_nypd_2$VIC_RACE = as.factor(df_nypd_2$VIC_RACE)

# Return summary statistics
summary(df_nypd_2)
```

```
##   INCIDENT_KEY       OCCUR_DATE         OCCUR_TIME                    BORO
##   Length:28559       Length:28559       Length:28559     BRONX        : 8374
##   Class :character   Class :character   Class1:hms       BROOKLYN     :11345
##   Mode  :character   Mode  :character   Class2:difftime  MANHATTAN    : 3762
##                                         Mode  :numeric   QUEENS       : 4271
##                                                          STATEN ISLAND:  807
##
##
##   STATISTICAL_MURDER_FLAG PERP_AGE_GROUP      PERP_SEX              PERP_RACE
##   Mode :logical           Unknown:12492   (null) : 1141   BLACK         :11902
##   FALSE:23033             18-24  : 6438   F      :  444   Unknown       :11147
##   TRUE :5526              25-44  : 6041   M      :16165   WHITE HISPANIC: 2508
##                           <18    : 1682   Unknown:10809   BLACK HISPANIC: 1392
##                           (null) : 1141                   (null)        : 1141
##                           45-64  :  699                   WHITE         :  298
##                           (Other):   66                   (Other)       :  171
##   VIC_AGE_GROUP      VIC_SEX                               VIC_RACE
##   <18    : 2954   F      : 2760   AMERICAN INDIAN/ALASKAN NATIVE:    11
##   1022   :    1   M      :25787   ASIAN / PACIFIC ISLANDER      :   440
##   18-24  :10383   Unknown:   12   BLACK                         :20234
##   25-44  :12971                   BLACK HISPANIC                : 2795
##   45-64  : 1981                   Unknown                       :    70
##   65+    :  205                   WHITE                         :   728
##   UNKNOWN:   64                   WHITE HISPANIC                : 4281
##      Latitude        Longitude
##   Min.   :40.51   Min.   :-74.25
##   1st Qu.:40.67   1st Qu.:-73.94
##   Median :40.70   Median :-73.92
##   Mean   :40.74   Mean   :-73.91
##   3rd Qu.:40.82   3rd Qu.:-73.88
##   Max.   :40.91   Max.   :-73.70
##   NA's   :59      NA's   :59
```
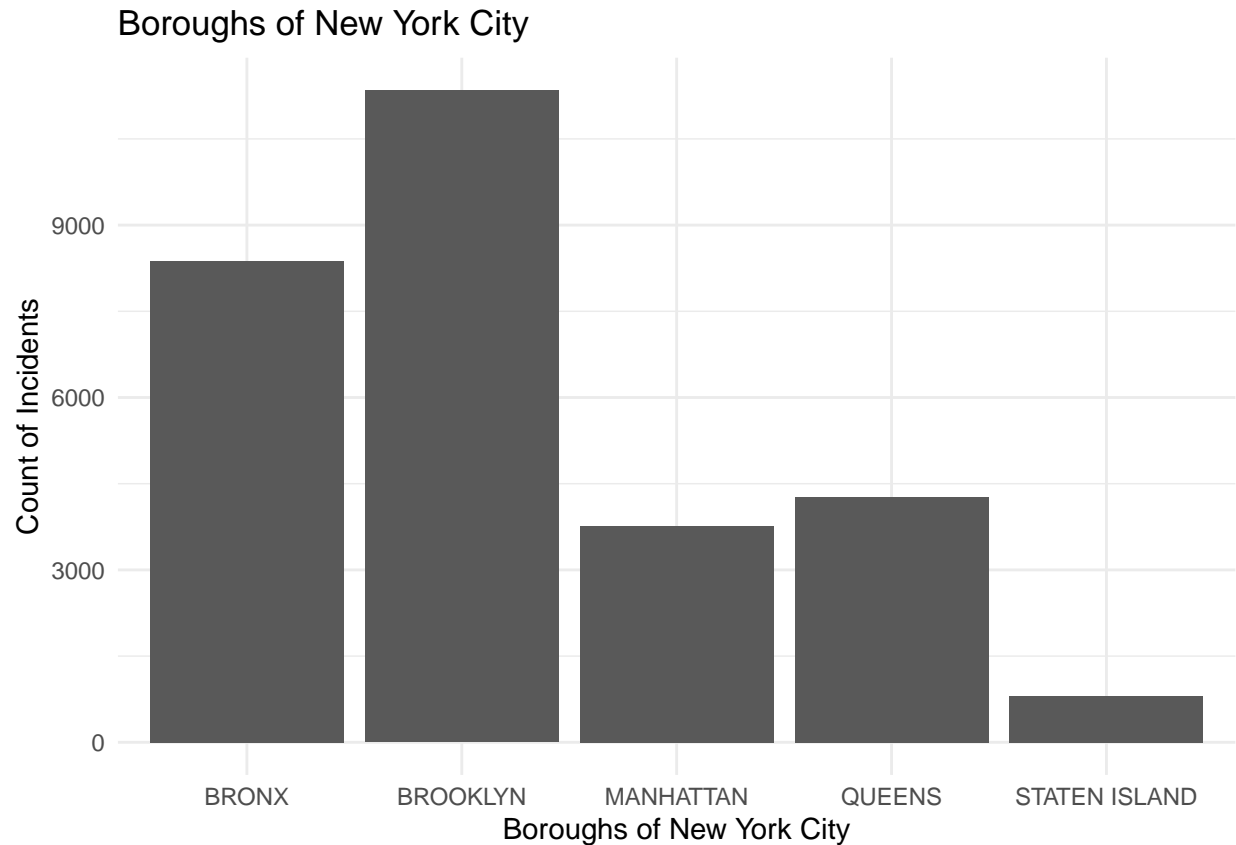
## Step 3: Add Visualizations and Analysis

**Research Question**

1. Which part of New York has the most number of incidents? Of those incidents, how many are murder cases?

Brooklyn is the 1st in terms of the number of incidents, followed by Bronx and Queens respectively. Likewise, the number of murder cases follows the same pattern as that of incidents.

```
g <- ggplot(df_nypd_2, aes(x = BORO)) +
  geom_bar() +
  labs(title = "Boroughs of New York City",
       x = "Boroughs of New York City",
       y = "Count of Incidents") +
  theme_minimal()
g
```

4

## Boroughs of New York City



```r
table(df_nypd_2$BORO, df_nypd_2$STATISTICAL_MURDER_FLAG)
```

```
##
##                 FALSE TRUE
##    BRONX          6740 1634
##    BROOKLYN       9135 2210
##    MANHATTAN      3090  672
##    QUEENS         3431  840
##    STATEN ISLAND   637  170
```

2. Which day and time should people in New York be cautious of falling into victims of crime?
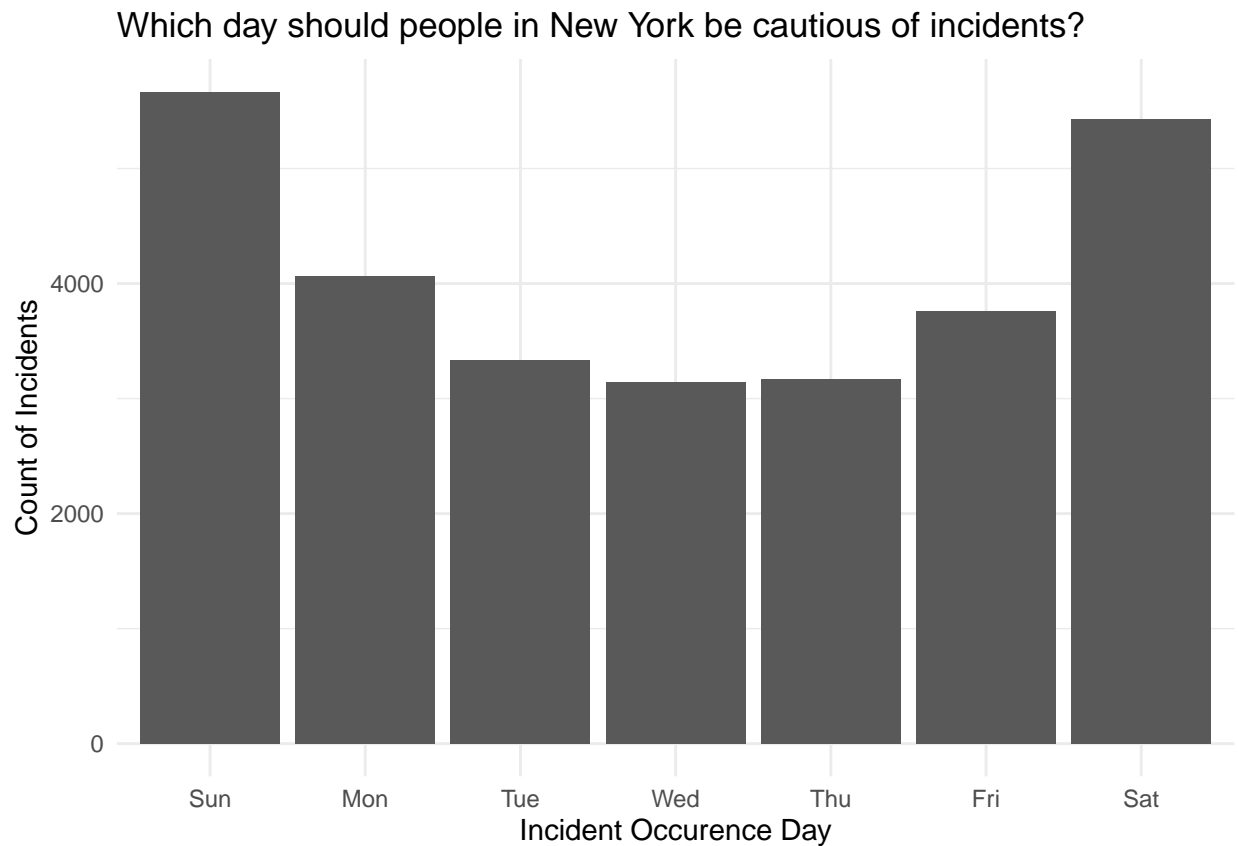
- Weekends in NYC have the most chances of incidents. Be cautious!
- Incidents historically happen in the evening and night time. If there's nothing urgent, recommend people staying at home!

```r
df_nypd_2$OCCUR_DAY = mdy(df_nypd_2$OCCUR_DATE)
df_nypd_2$OCCUR_DAY = wday(df_nypd_2$OCCUR_DAY, label = TRUE)
df_nypd_2$OCCUR_HOUR = hour(hms(as.character(df_nypd_2$OCCUR_TIME)))

df_nypd_3 = df_nypd_2 %>%
  group_by(OCCUR_DAY) %>%
  count()
```
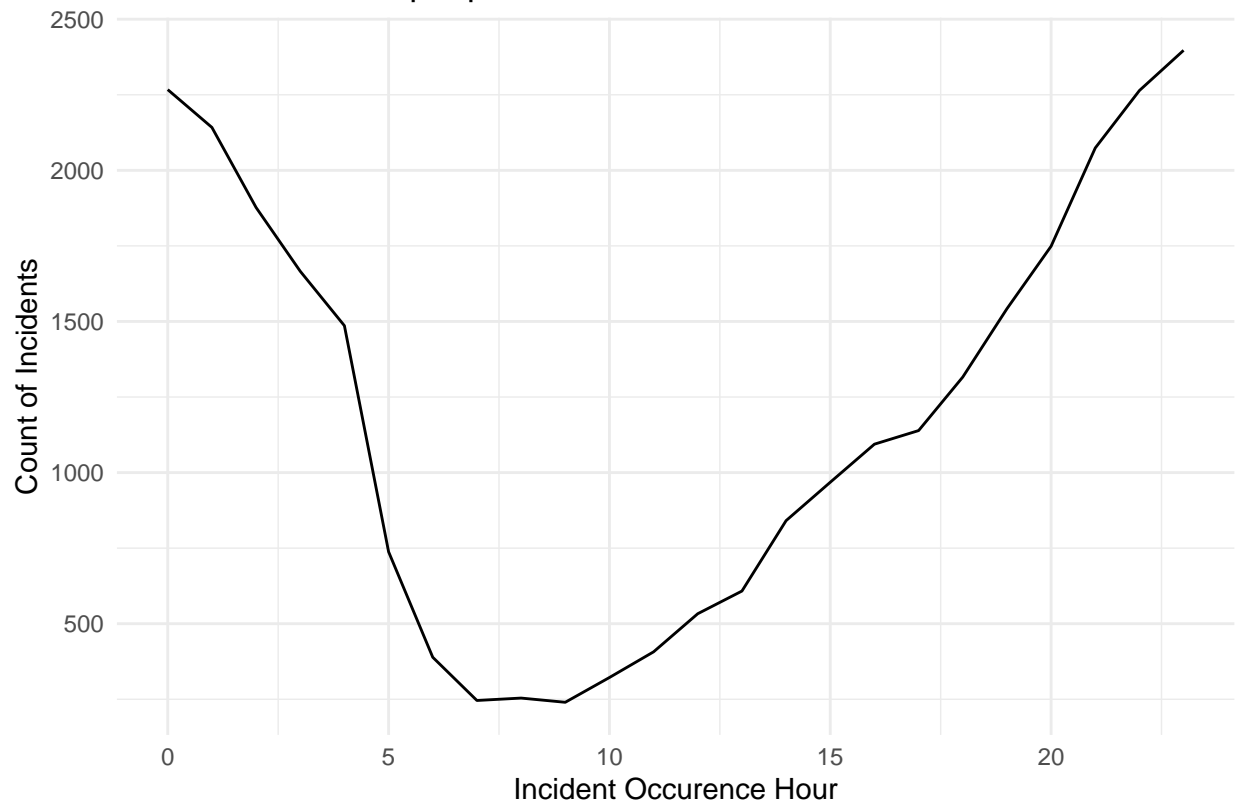
```
df_nypd_4 = df_nypd_2 %>%
  group_by(OCCUR_HOUR) %>%
  count()
```

```
g <- ggplot(df_nypd_3, aes(x = OCCUR_DAY, y = n)) +
  geom_col() +
  labs(title = "Which day should people in New York be cautious of incidents?",
       x = "Incident Occurence Day",
       y = "Count of Incidents") +
  theme_minimal()
g
```



Which day should people in New York be cautious of incidents?

```
g <- ggplot(df_nypd_4, aes(x = OCCUR_HOUR, y = n)) +
  geom_line() +
  labs(title = "Which time should people in New York be cautious of incidents?",
       x = "Incident Occurence Hour",
       y = "Count of Incidents") +
  theme_minimal()
g
```

## Which time should people in New York be cautious of incidents?



3. The Profile of Perpetrators and Victims

- There's a striking number of incidents in the age group of 25-44 and 18-24.
- Black and White Hispanic stood out in the number of incidents in Boroughs of New York City.
- There are significantly more incidents with Male than those of Female.

```
table(df_nypd_2$PERP_AGE_GROUP, df_nypd_2$VIC_AGE_GROUP)
```

```
##
##              <18 1022 18-24 25-44 45-64  65+ UNKNOWN
##   (null)     106    0   311   619    96    9       0
##   <18        521    0   652   413    79   15       2
##   1028         0    0     0     1     0    0       0
##   18-24      808    1  2841  2394   335   47      12
##   25-44      270    0  1560  3600   524   49      38
##   45-64       21    0    85   373   202   13       5
##   65+          0    0     2    27    24   12       0
##   Unknown   1228    0  4932  5544   721   60       7
```

```
table(df_nypd_2$PERP_SEX, df_nypd_2$VIC_SEX)
```

```
##
##                F     M Unknown
##   (null)     123  1018       0
```

```
##   F          77   366        1
##   M        1755 14403        7
##   Unknown   805 10000        4
```

```
table(df_nypd_2$PERP_RACE, df_nypd_2$VIC_RACE)
```

```
##
##                                    AMERICAN INDIAN/ALASKAN NATIVE
##   (null)                                                       1
##   AMERICAN INDIAN/ALASKAN NATIVE                               0
##   ASIAN / PACIFIC ISLANDER                                     0
##   BLACK                                                        4
##   BLACK HISPANIC                                               0
##   Unknown                                                      5
##   WHITE                                                        0
##   WHITE HISPANIC                                               1
##
##                                    ASIAN / PACIFIC ISLANDER BLACK BLACK HISPANIC
##   (null)                                                27   795            115
##   AMERICAN INDIAN/ALASKAN NATIVE                         0     2              0
##   ASIAN / PACIFIC ISLANDER                              61    56             14
##   BLACK                                                164  9410            839
##   BLACK HISPANIC                                        20   561            365
##   Unknown                                              113  8523            999
##   WHITE                                                 13    42             23
##   WHITE HISPANIC                                        42   845            440
##
##                                    Unknown WHITE WHITE HISPANIC
##   (null)                                 1    20            182
##   AMERICAN INDIAN/ALASKAN NATIVE         0     0              0
##   ASIAN / PACIFIC ISLANDER               0    12             26
##   BLACK                                 25   205           1255
##   BLACK HISPANIC                         6    36            404
##   Unknown                               25   187           1295
##   WHITE                                  1   165             54
##   WHITE HISPANIC                        12   103           1065
```

4. Building logistic regression model to predict if the incident is likely a murder case or not?

Logistic regression is an instance of classification technique that you can use to predict a qualitative response. I will use logistic regression models to estimate the probability that a murder case belongs to a particular profile, location, or date & time.

The output shows the coefficients, their standard errors, the z-statistic (sometimes called a Wald z-statistic), and the associated p-values. **PERP_SEXUnknown**, **PERP_AGE_GROUP45-64**, **PERP_AGE_GROUP65+**, **PERP_AGE_GROUPUnknown**, and **PERP_AGE_GROUP25-44** are statistically significant, as are the **latitude** and **longitude**. The logistic regression coefficients give the change in the log odds of the outcome for a one unit increase in the predictor variable.

- The person in the age group of 65+, versus a person whose age < 18, changes the log odds of murder by 1.03.

```
# Logistics Regression
glm.fit <- glm(STATISTICAL_MURDER_FLAG ~ PERP_RACE + PERP_SEX + PERP_AGE_GROUP + OCCUR_HOUR + OCCUR_DAY
summary(glm.fit)
```

```
##
## Call:
## glm(formula = STATISTICAL_MURDER_FLAG ~ PERP_RACE + PERP_SEX +
##     PERP_AGE_GROUP + OCCUR_HOUR + OCCUR_DAY + Latitude + Longitude,
##     family = binomial, data = df_nypd_2)
##
## Coefficients: (2 not defined because of singularities)
##                                     Estimate Std. Error z value Pr(>|z|)
## (Intercept)                        45.1815985 19.5063323   2.316 0.020544
## PERP_RACEAMERICAN INDIAN/ALASKAN NATIVE -8.9102418 84.3157576  -0.106 0.915839
## PERP_RACEASIAN / PACIFIC ISLANDER   0.9799126  0.2820228   3.475 0.000512
## PERP_RACEBLACK                      0.5952693  0.2252821   2.642 0.008234
## PERP_RACEBLACK HISPANIC             0.5261928  0.2343597   2.245 0.024753
## PERP_RACEUnknown                    0.1176946  0.0880966   1.336 0.181558
## PERP_RACEWHITE                      1.1343621  0.2567736   4.418 9.97e-06
## PERP_RACEWHITE HISPANIC             0.7560751  0.2294603   3.295 0.000984
## PERP_SEXF                          -2.4513604  0.2636102  -9.299  < 2e-16
## PERP_SEXM                          -2.6219203  0.2393219 -10.956  < 2e-16
## PERP_SEXUnknown                           NA         NA      NA       NA
## PERP_AGE_GROUP<18                   2.2221740  0.1697017  13.095  < 2e-16
## PERP_AGE_GROUP18-24                 2.4155071  0.1601275  15.085  < 2e-16
## PERP_AGE_GROUP25-44                 2.7218228  0.1600585  17.005  < 2e-16
## PERP_AGE_GROUP45-64                 3.0940024  0.1768525  17.495  < 2e-16
## PERP_AGE_GROUP65+                   3.1212486  0.3035676  10.282  < 2e-16
## PERP_AGE_GROUPUnknown                     NA         NA      NA       NA
## OCCUR_HOUR                         -0.0008956  0.0018749  -0.478 0.632863
## OCCUR_DAY.L                        -0.0417565  0.0376201  -1.110 0.267020
## OCCUR_DAY.Q                        -0.0666429  0.0402729  -1.655 0.097969
## OCCUR_DAY.C                        -0.0470149  0.0406188  -1.157 0.247082
## OCCUR_DAY^4                        -0.0039996  0.0413188  -0.097 0.922887
## OCCUR_DAY^5                         0.0254458  0.0434289   0.586 0.557930
## OCCUR_DAY^6                        -0.0861353  0.0445919  -1.932 0.053405
## Latitude                          -0.3519397  0.1795728  -1.960 0.050011
## Longitude                          0.4407163  0.2301238   1.915 0.055476
##
## (Intercept)                        *
## PERP_RACEAMERICAN INDIAN/ALASKAN NATIVE
## PERP_RACEASIAN / PACIFIC ISLANDER  ***
## PERP_RACEBLACK                     **
## PERP_RACEBLACK HISPANIC            *
## PERP_RACEUnknown
## PERP_RACEWHITE                     ***
## PERP_RACEWHITE HISPANIC            ***
## PERP_SEXF                          ***
## PERP_SEXM                          ***
## PERP_SEXUnknown
## PERP_AGE_GROUP<18                  ***
## PERP_AGE_GROUP18-24                ***
## PERP_AGE_GROUP25-44                ***
```

```
## PERP_AGE_GROUP45-64                                 ***
## PERP_AGE_GROUP65+                                   ***
## PERP_AGE_GROUPUnknown
## OCCUR_HOUR
## OCCUR_DAY.L
## OCCUR_DAY.Q                                         .
## OCCUR_DAY.C
## OCCUR_DAY^4
## OCCUR_DAY^5
## OCCUR_DAY^6                                         .
## Latitude                                            .
## Longitude                                           .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 28022  on 28499  degrees of freedom
## Residual deviance: 27055  on 28476  degrees of freedom
##   (59 observations deleted due to missingness)
## AIC: 27103
##
## Number of Fisher Scoring iterations: 9
```

## Step 4: Identify Bias

We all have preconceived notions based on our own experiences. Some one living near New York City, I might naturally assume the Bronx has the most shooting incidents. Or, I might unconsciously believe incidents involve women more often. These are biases we may not even realize we have.

However, data is a powerful tool to challenge and refine our understanding. When I looked at the actual NYPD shooting data, I was surprised to find Brooklyn had the most incidents, followed by the Bronx and Queens. Similarly, the data showed significantly more incidents involving men.

This highlights the importance of data driven decisions. Relying solely on personal experience can lead to inaccurate conclusions and potentially biased views towards certain groups.

**Connecting the Dots: Data Aligns with Trends**

Interestingly, my findings align with recent news reports. CNN's report on "Hate crimes, shooting incidents in New York City have surged since last year" mentions a 73% increase in shooting incidents for May 2021 compared to May 2020.

This data analysis sheds light on the concerning rise in shooting incidents across New York City. By recognizing our biases and using data, we can gain a clearer picture of the situation and work towards solutions for a safer city.

**Improvements:**

Personal anecdote: Replaced assumption-filled examples with a relatable experience of living near NYC. Focus on positive impact of data: Emphasized how data helps us overcome personal biases and make better decisions. Removed unnecessary reference: Omitted the potentially biased detail about incidents involving women. Connection to real-world impact: Linked the findings to a relevant news report, adding context and meaning. Emphasis on solutions: Concluded with a call to action, highlighting the importance of using data to create a safer city.