# NYPD Shooting Incident

*TSP*

*10/29/2022*

#Peer-graded Assignement: NYPD Shooting Incident Data Report

Assignement Tasks: Import, tidy and analyze the NYPD Shooting Incident dataset obtained. Be sure your project is reproducible and contains some visualization and analysis.You may use the data to do any analysis that is of interest to you. You should include at least two visualizations and one model. Be sure to identify any bias possible in the data and in your analysis.

#Step 1: Install packages and enable the package required for tjis analysis

#Step 2: Gether NYPD Shooting incident data in the csv. format from the URL The data is a plubic data which is avialble in Data Gov. To access the data you can go to https://catalog.data.gov/dataset and find the dataset titled NYPD Shooting Incident Data (Historic). Here I have copied the URL link of the data in csv format and you r command to download it for futher analysis.

#Step 3: Perform basic analysis and data cleansing

#Step 4: Subsetting and Build a clustering model I will use only lat and lon columns to perfrom geospatial clustering based on location (lat,lon). So, the first step I will subsetting data and then perfrom Basic K-mean clustering

```r
# Subsetting the NYPD data to be only lat and lon
df <- NYPD_shooting_incident[,17:18]

# Set the seed for reproducable
set.seed(123)
#Find the optimal k with silhouette method
fviz_nbclust(df, kmeans, method = "silhouette")
```

## Optimal number of clusters



```
#Perfrom kmean based on kmean clustering
final <- kmeans(df, 3)
print(final)
```

```
## K-means clustering with 3 clusters of sizes 887, 3109, 3247
##
## Cluster means:
##   Latitude Longitude
## 1 40.67004 -73.79776
## 2 40.82825 -73.90669
## 3 40.66741 -73.95509
##
## Clustering vector:
##    [1] 2 2 1 2 2 2 2 2 2 3 3 2 2 2 3 3 3 3 2 2 3 2 2 3 2 2 1 1 3 2 2 2 2 2
##   [35] 1 1 2 3 3 2 2 2 2 2 2 2 2 2 2 2 2 3 2 1 2 3 3 2 2 2 3 2 2 3 2 2 3 3 2
##   [69] 1 3 3 1 3 2 2 1 1 3 2 1 2 3 1 1 3 2 2 3 3 3 2 3 3 2 2 2 2 2 3 1 3 3
##  [103] 2 3 2 3 2 3 2 3 1 2 1 2 3 2 2 1 3 2 2 1 2 2 2 3 3 2 2 2 2 2 2 2 2 2
##  [137] 1 2 3 1 2 1 3 2 3 2 2 3 3 2 3 1 2 1 1 2 2 2 3 2 3 3 2 2 3 1 2 1 2 3
##  [171] 2 2 3 2 3 2 1 2 2 2 2 3 2 3 3 2 2 2 2 3 3 1 2 2 2 3 2 2 3 3 2 3 2 3
##  [205] 3 1 1 3 2 3 2 2 2 1 2 3 3 3 3 2 2 3 2 2 3 2 2 2 2 2 2 1 3 2 2 3 2 2 3
##  [239] 3 1 3 1 1 2 2 1 3 2 3 2 2 2 2 2 3 2 1 2 2 3 2 1 2 2 2 3 3 2 2 2 3 1
##  [273] 3 3 2 3 2 3 3 2 2 3 1 3 3 2 2 2 2 3 1 2 3 2 3 2 3 3 2 1 2 2 1 3 2 3
##  [307] 2 2 3 2 2 1 2 3 3 3 2 2 3 1 1 3 3 2 2 1 3 2 1 3 3 2 1 2 2 3 2 1 2 2
##  [341] 3 3 3 3 3 2 3 3 2 2 2 2 2 3 3 3 2 2 2 2 3 3 3 3 3 2 1 3 1 1 3 2 3 2 3
##  [375] 3 1 3 2 3 3 2 1 2 3 3 2 2 1 3 3 2 3 3 2 1 3 2 3 2 3 2 2 3 3 2 2 2 3
##  [409] 1 2 2 2 3 3 1 3 2 2 1 3 1 2 3 3 1 2 2 3 3 2 2 2 2 3 3 3 1 2 2 3 3 3
```

```
##  [443] 3 2 2 2 2 2 1 1 1 3 2 2 2 2 1 2 3 3 2 2 1 1 2 3 1 2 2 3 2 3 2 2 3 3
##  [477] 1 2 3 3 2 3 2 2 1 2 3 3 3 3 2 1 3 3 3 3 2 3 1 2 3 3 3 1 3 3 2 3 2 3 2
##  [511] 3 2 3 1 3 2 3 3 2 3 3 3 1 1 2 2 2 2 3 3 3 3 2 3 2 2 2 2 2 3 3 2 2 2
##  [545] 2 2 3 2 2 1 3 2 2 3 3 2 3 3 3 2 3 3 2 2 2 3 3 3 3 3 2 3 1 3 2 2 3 2
##  [579] 3 3 2 3 3 3 1 3 3 3 3 2 3 2 1 3 3 2 3 1 1 3 3 3 2 1 2 2 3 3 3 1 2 1
##  [613] 3 2 3 3 3 1 3 3 2 3 2 2 3 3 1 3 1 1 3 1 2 3 3 3 2 3 3 3 2 3 3 3 2 2
##  [647] 2 2 2 2 3 1 1 2 2 3 2 3 3 2 3 3 2 2 1 2 3 2 2 3 3 3 2 2 2 2 3 2 3 2
##  [681] 1 3 3 1 1 2 2 3 3 2 1 2 3 3 2 2 2 2 2 3 3 3 3 3 3 3 3 2 2 3 3 3 1 2
##  [715] 2 3 2 3 3 3 2 3 2 3 2 3 3 2 3 3 1 2 3 2 3 3 2 3 2 3 2 3 2 3 3 2 2 3 3 2
##  [749] 2 2 3 2 2 2 3 3 2 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 3 2 2 1 3 3 2 3 3 3 3
##  [783] 3 2 1 3 3 3 3 3 2 3 3 1 2 2 3 3 2 2 2 3 3 2 2 2 3 3 2 2 2 3 3 2 2 2
##  [817] 2 3 2 2 2 2 2 2 3 3 1 2 2 2 2 3 3 3 1 1 2 3 3 2 1 2 1 2 2 2 3 2 2 2
##  [851] 2 2 1 3 3 3 2 3 2 2 3 1 1 1 3 3 2 3 2 3 3 3 3 3 2 2 1 3 2 2 3 3 2 2
##  [885] 2 2 1 2 3 3 3 3 3 3 3 3 1 3 1 3 1 3 3 2 3 2 1 2 2 3 3 1 3 3 2 3 3 2
##  [919] 2 2 3 3 2 2 3 1 3 2 2 3 3 3 2 2 2 3 1 2 3 2 2 2 2 3 3 3 1 3 1 2 3 2
##  [953] 2 3 3 3 3 1 3 2 3 2 3 2 2 3 3 3 2 3 3 3 2 3 3 2 1 2 1 2 2 2 2 2 2 3
##  [987] 2 2 3 2 2 2 3 3 1 3 2 1 3 3 3 2 3 3 3 2 3 3 2 2 3 3 3 3 2 3 3 2 2 3
## [1021] 1 3 2 2 2 3 3 3 3 1 2 1 2 2 2 3 2 2 2 3 2 3 2 3 2 2 2 2 2 1 1 3 2 2
## [1055] 2 3 3 2 2 3 3 2 3 2 3 3 2 3 1 3 2 3 2 3 2 3 2 3 1 3 2 3 2 2 3 3 1 3
## [1089] 1 2 2 2 2 2 3 2 2 2 3 3 3 3 3 2 2 2 3 3 3 2 3 2 2 1 3 2 3 2 2 2 1 2
## [1123] 3 2 2 3 3 3 1 2 3 2 2 2 3 2 3 3 2 3 2 3 2 3 1 2 2 3 3 1 2 2 3 3 3 2
## [1157] 3 3 3 3 3 2 2 2 3 3 1 2 2 3 3 2 2 1 2 3 3 1 3 2 3 3 3 1 2 1 3 3 3 3
## [1191] 2 3 3 2 2 2 2 2 2 2 2 2 2 3 1 2 3 3 3 3 3 1 3 3 2 3 1 2 3 2 2 2 2 3 2
## [1225] 3 2 3 3 2 2 2 2 2 3 3 2 2 2 3 1 3 2 2 3 3 3 3 3 3 3 1 2 2 2 3 2 2 1 3
## [1259] 3 2 2 2 2 3 3 2 3 1 3 2 3 2 2 1 1 2 2 3 2 3 1 2 2 3 3 2 3 3 3 1 3 2
## [1293] 2 2 3 3 2 3 3 2 2 2 2 1 3 3 2 3 3 3 2 1 3 1 2 3 2 3 2 3 2 2 1 2 1 3
## [1327] 3 2 2 2 3 2 3 2 3 3 3 2 2 1 3 2 2 2 2 3 3 2 2 1 2 2 3 3 2 3 2 1 3 3
## [1361] 3 3 3 3 3 3 2 3 2 1 3 3 3 3 3 3 1 3 2 3 3 3 1 3 3 2 2 3 3 3 3 2 2 3
## [1395] 2 3 3 3 3 2 2 3 3 2 3 3 2 1 2 1 3 2 3 3 2 2 2 2 3 3 2 2 3 1 2 1 3 3
## [1429] 2 3 2 1 3 3 2 2 2 3 1 1 2 3 3 2 3 3 3 3 2 2 3 3 1 3 2 3 3 1 3 3 2 1
## [1463] 2 2 2 2 2 3 1 3 3 3 3 3 2 3 2 2 3 2 2 2 1 2 1 2 3 1 1 3 2 1 3 2 2 3
## [1497] 2 3 2 2 3 3 3 3 1 3 1 3 2 3 2 3 2 2 3 2 3 1 2 2 2 3 3 2 1 3 3 3 1 2
## [1531] 3 2 3 3 3 2 2 3 2 2 1 3 2 2 2 3 2 2 3 2 3 2 3 3 3 3 2 2 2 2 3 2 2 2
## [1565] 1 1 1 2 3 3 1 2 2 3 2 1 1 3 1 3 2 1 1 1 2 2 2 3 1 3 3 1 2 2 2 2 2 1
## [1599] 3 3 2 1 2 2 3 2 3 3 2 2 2 3 2 3 1 3 3 3 2 1 2 3 1 3 2 3 2 1 3 1 3 3
## [1633] 2 1 3 3 2 2 2 2 3 1 3 3 3 1 2 3 3 3 2 3 2 3 2 3 3 2 2 2 3 3 2 3 2 1
## [1667] 3 3 2 3 2 3 3 2 1 3 3 2 3 2 3 2 3 2 2 2 3 3 3 2 3 1 2 3 2 3 3 2 3 2 1
## [1701] 3 1 2 2 2 2 2 1 1 3 2 3 2 2 2 2 1 2 2 3 3 3 2 3 3 2 2 3 3 3 3 2 3 2
## [1735] 2 2 3 1 3 3 3 3 2 1 2 2 2 2 3 3 2 2 2 2 3 3 3 3 2 3 2 3 3 3 3 3 2 3
## [1769] 1 3 2 3 2 2 2 2 3 3 2 3 3 3 3 1 3 2 1 3 3 2 3 3 3 2 2 2 2 3 2 3 1 2
## [1803] 2 2 2 2 2 2 3 1 3 3 3 3 3 2 2 2 2 3 2 2 3 2 2 3 3 3 2 2 3 2 2 2 3 3
## [1837] 1 2 3 3 3 2 2 3 2 2 1 2 1 2 2 2 2 3 2 1 3 3 2 3 3 1 2 3 1 3 2 1 1 3
## [1871] 3 2 3 2 2 2 3 2 3 2 3 3 3 3 1 3 3 2 2 2 2 3 2 3 3 3 2 2 3 3 2 1 1 3
## [1905] 3 2 3 3 2 2 2 1 2 2 2 3 1 3 3 2 2 2 2 3 2 2 1 3 2 2 3 3 3 2 3 1 3 2
## [1939] 3 2 3 2 3 3 1 3 3 2 3 2 2 3 2 3 3 3 2 3 2 3 2 2 3 2 2 2 1 1 2 2 3 3
## [1973] 2 3 2 2 1 3 2 2 2 3 3 3 3 3 3 2 3 2 3 2 2 3 2 1 3 3 2 2 1 3 3 2 3 3
## [2007] 3 3 1 1 1 2 3 3 3 1 2 3 2 2 2 2 1 2 3 3 1 2 3 3 3 1 3 2 2 2 2 3 3 2
## [2041] 2 3 3 3 3 2 2 3 3 2 3 3 3 3 2 3 2 2 1 2 2 3 1 3 2 3 2 3 1 2 2 2 1 2
## [2075] 3 3 2 2 2 1 3 3 2 2 1 2 3 2 2 3 3 3 1 3 1 3 3 3 3 1 3 2 3 1 3 3 2 3
## [2109] 2 3 3 2 3 3 2 3 2 2 3 3 2 3 3 2 3 2 3 3 3 2 2 3 1 2 3 3 2 2 3 3 2 3
## [2143] 2 2 3 2 2 1 3 3 2 3 2 3 2 2 2 3 3 3 3 3 1 1 1 3 3 2 2 1 2 2 3 1 3 3
## [2177] 2 3 3 2 2 3 3 2 2 2 3 2 1 2 3 3 2 2 3 1 2 2 2 2 2 2 3 3 3 2 1 2 2 2
## [2211] 1 2 2 3 2 3 2 2 1 3 3 1 3 2 3 2 3 3 2 2 3 3 2 3 1 1 1 1 2 2 3 3 3 2
## [2245] 3 2 2 1 3 3 2 3 3 2 3 2 3 2 3 3 2 2 2 3 3 3 2 1 3 3 3 2 2 2 3 3 2 3
```

3

```
## [2279] 3 1 3 2 2 1 2 3 3 2 3 1 1 3 3 3 2 3 2 3 3 2 3 3 1 1 3 2 2 3 3 2 1 3
## [2313] 3 3 2 3 2 3 2 3 3 3 2 3 2 3 3 3 3 2 2 3 2 3 1 3 2 2 3 3 2 1 2 3 1 1
## [2347] 3 3 2 3 3 3 3 3 3 1 3 2 2 2 1 3 2 3 2 2 2 2 3 3 3 3 3 2 1 3 3 2 3 2
## [2381] 2 2 1 1 2 3 2 2 3 2 2 1 3 3 1 2 3 3 3 1 3 2 3 2 3 3 3 3 2 1 3 3 2 1
## [2415] 2 2 1 2 2 1 2 3 2 3 2 2 2 1 2 2 2 3 3 2 3 1 2 2 2 3 1 3 2 2 3 1 2 3
## [2449] 2 3 1 1 3 2 3 3 3 2 2 2 3 1 3 3 2 2 2 2 1 2 2 2 3 3 3 1 2 3 2 2 2 1
## [2483] 1 2 3 3 3 3 2 2 3 3 2 3 2 3 3 3 3 2 3 2 1 3 3 3 1 2 3 3 3 1 2 3 3 1
## [2517] 3 3 3 3 3 3 3 2 3 3 3 3 1 3 3 3 2 3 3 2 2 2 2 2 2 3 2 3 2 3 3 3 3
## [2551] 2 3 3 3 2 2 1 3 3 3 2 3 1 3 3 2 3 1 3 3 1 2 3 3 1 2 2 2 3 3 2 2 2 2
## [2585] 3 3 3 3 3 2 3 3 2 3 3 1 3 3 2 3 3 3 3 2 2 3 3 2 3 2 3 3 2 2 2 3 1 3
## [2619] 2 1 2 3 3 1 3 2 3 1 2 2 2 2 3 2 2 3 2 2 3 3 3 2 3 3 2 1 3 1 3 3 2 2
## [2653] 2 2 3 2 3 3 2 2 3 3 2 3 3 2 3 3 3 3 3 3 3 1 2 2 2 3 1 2 2 2 2 3 3 1
## [2687] 2 3 1 2 2 3 2 3 2 2 1 2 2 3 2 3 2 2 2 3 3 3 3 2 2 2 3 2 3 3 2 2 2 2
## [2721] 2 3 3 3 2 3 2 2 1 3 3 2 1 2 3 3 3 3 3 3 2 2 3 2 3 2 3 1 3 2 1 3 2 3
## [2755] 3 3 3 3 1 2 3 2 3 2 2 3 1 3 3 1 2 2 1 2 2 2 3 2 1 2 3 2 3 2 3 3 2 2
## [2789] 3 2 2 2 3 1 2 2 2 3 2 3 1 2 2 1 2 3 3 2 3 2 3 2 2 3 2 3 3 2 3 3 3 3
## [2823] 2 1 3 3 2 3 3 2 2 2 3 3 1 2 2 3 2 3 1 1 3 2 3 1 2 3 1 2 3 3 2 1 2 2
## [2857] 2 3 3 2 2 3 2 3 2 2 3 2 2 3 2 3 2 2 1 1 2 3 3 3 3 1 1 3 2 3 2 2 2 3 2
## [2891] 2 2 2 2 3 2 2 2 3 2 1 3 2 2 2 2 3 2 1 3 1 3 2 3 3 3 3 3 3 2 2 3 3 3
## [2925] 3 2 3 2 2 3 2 2 3 3 3 2 3 2 3 1 1 2 2 3 3 3 2 3 3 3 3 3 3 2 2 2 2 2
## [2959] 2 2 2 2 2 2 1 3 3 2 3 3 3 2 3 3 3 2 2 3 2 2 2 2 1 2 3 3 3 3 3 2 3 2
## [2993] 2 2 3 2 2 2 3 3 3 2 3 3 2 3 3 2 2 2 3 3 2 2 3 2 3 2 2 2 2 3 3 3 3 2
## [3027] 2 2 3 3 1 2 3 3 2 2 2 3 3 2 3 2 2 3 3 3 3 2 2 2 2 1 2 2 2 1 2 1 3 3
## [3061] 2 3 1 2 3 3 3 2 3 1 2 2 2 2 3 3 3 3 2 1 2 3 1 2 2 3 2 3 1 3 2 1 2 2 2
## [3095] 1 3 3 3 3 3 2 2 3 1 3 2 2 3 3 2 3 2 3 2 3 2 3 3 1 3 3 1 3 2 2 2 3 3
## [3129] 3 2 2 3 3 3 2 2 3 3 2 3 2 3 2 3 2 3 2 1 3 2 1 2 2 2 3 3 3 2 3 3 3 3
## [3163] 1 2 1 2 3 1 2 2 3 2 2 2 3 2 2 3 2 3 3 3 2 3 3 2 3 3 3 1 2 3 2 3 2 2
## [3197] 2 3 2 3 2 3 3 2 3 2 2 3 2 2 3 3 2 2 1 2 3 2 3 3 2 2 3 3 3 2 2 2 2 2
## [3231] 3 2 2 1 3 2 1 2 2 2 3 3 2 1 1 3 3 3 3 3 3 3 1 3 1 3 3 1 2 3 1 2 1 1
## [3265] 2 3 2 2 2 3 2 2 3 2 3 2 1 2 2 2 2 3 3 2 1 1 1 1 3 2 2 2 3 3 3 3 3 2
## [3299] 3 1 3 1 2 2 3 2 3 3 3 3 2 2 1 3 2 2 3 2 2 2 3 3 2 2 2 2 3 2 2 3 3 3
## [3333] 1 3 2 3 2 3 1 3 2 2 2 2 3 2 3 2 2 3 1 2 3 2 2 2 3 3 3 3 2 3 3 2 3 3
## [3367] 1 3 1 3 1 1 2 3 2 3 2 3 2 2 3 3 2 2 1 1 2 3 1 3 3 3 2 1 3 3 3 2 3 2
## [3401] 1 3 1 2 3 3 1 3 2 3 1 2 3 2 2 2 2 2 3 3 3 2 3 3 3 3 3 2 3 1 2 3 1 2
## [3435] 3 3 1 2 2 1 1 2 2 3 3 1 3 2 1 2 1 3 2 2 3 3 3 3 3 3 1 2 2 2 2 3 2 3
## [3469] 3 3 3 2 1 3 1 2 3 1 3 3 3 3 3 3 2 3 2 1 3 1 3 3 1 2 2 1 2 2 1 2 3 1
## [3503] 1 2 1 2 3 3 2 2 2 3 3 1 3 3 2 3 2 2 2 1 2 3 3 3 1 2 3 3 2 3 2 3 2 3
## [3537] 3 2 3 2 1 3 3 2 3 2 2 2 3 2 2 2 2 1 3 2 3 2 3 2 2 3 3 3 2 2 2 3 2 3
## [3571] 3 3 3 1 3 2 3 2 1 2 3 2 3 1 3 1 3 3 3 2 2 3 2 2 2 3 2 2 2 3 3 2 3 1
## [3605] 3 2 3 3 3 3 3 3 3 2 3 3 3 3 3 2 3 1 3 2 2 2 3 3 3 3 3 3 3 2 2 3 3 2
## [3639] 3 1 2 3 3 2 3 3 1 3 1 2 3 3 3 3 2 2 3 1 3 3 3 2 2 2 3 3 2 2 2 2 2 2
## [3673] 3 3 2 3 2 1 3 3 3 3 2 3 3 3 2 3 2 3 3 2 2 1 2 3 3 3 3 1 2 3 2 2 2 3
## [3707] 3 2 2 1 1 3 2 2 1 3 3 2 2 2 1 2 2 1 2 3 2 2 3 3 1 2 3 2 2 3 1 3 2 2
## [3741] 3 2 2 3 2 2 3 3 2 2 3 2 1 3 1 2 2 2 3 3 2 2 3 3 1 2 2 3 2 1 3 1 1 1
## [3775] 2 1 3 3 1 3 3 2 3 3 2 3 3 3 2 2 3 2 2 3 3 3 3 2 2 2 3 2 2 2 2 2 1 2
## [3809] 2 3 3 2 3 2 2 3 3 3 3 3 3 3 2 2 2 2 2 2 2 3 3 3 2 3 3 3 3 3 2 2 2 2
## [3843] 1 3 3 3 1 2 3 2 3 1 1 1 1 3 2 3 3 3 3 1 3 3 3 3 2 1 2 3 3 2 1 3 3 3
## [3877] 3 2 3 2 2 3 1 1 3 3 2 2 3 3 2 2 2 2 3 3 3 3 2 1 2 3 2 3 2 3 2 2 3 3 3
## [3911] 1 2 1 3 2 2 3 2 2 2 3 2 2 2 3 1 2 3 3 3 1 3 2 3 2 3 1 3 2 3 2 2 2 3
## [3945] 3 3 2 2 2 2 2 3 3 3 2 3 2 2 1 3 3 2 3 3 2 2 3 1 3 2 2 1 2 3 3 2 2
## [3979] 2 1 3 3 2 2 1 2 3 2 3 3 3 3 2 2 2 2 1 3 3 3 1 2 1 3 3 2 3 2 2 2 3 1
## [4013] 2 1 2 2 3 3 1 3 3 2 2 3 2 3 3 2 3 1 3 2 2 2 1 1 2 2 3 2 3 2 1 3 3 3
## [4047] 2 2 3 3 2 3 3 2 3 1 2 3 2 3 2 2 2 3 3 2 2 3 2 3 3 1 1 1 1 3 2 2 2 3 3
## [4081] 2 2 2 3 3 2 3 3 2 2 2 3 2 3 2 3 2 2 3 3 3 3 3 3 3 3 2 2 3 3 3 3 3 2
```

4

```
## [4115] 1 3 3 2 2 3 3 2 3 2 1 3 3 2 3 1 1 3 1 3 2 3 2 1 3 2 3 2 2 2 2 1 2 3
## [4149] 3 2 3 2 2 3 2 3 2 2 3 1 3 2 3 3 3 3 3 2 3 2 3 3 3 3 2 3 2 3 3 3 3 2
## [4183] 3 2 2 2 3 3 2 3 2 1 1 2 3 2 2 2 2 2 3 3 1 2 3 3 2 3 1 1 2 2 1 3 3 2
## [4217] 2 2 1 3 3 3 3 3 3 1 2 3 2 2 2 2 2 2 2 3 3 3 3 3 3 3 2 3 2 1 3 3 1
## [4251] 2 3 3 2 2 2 3 2 3 3 2 2 3 3 3 3 3 3 3 3 3 3 3 3 2 3 2 2 2 3 2 1 3 3
## [4285] 2 2 1 2 3 2 1 1 2 2 2 2 2 1 3 2 3 1 2 3 2 3 2 2 3 2 2 2 2 3 2 2 3 3
## [4319] 3 2 2 3 3 2 1 2 2 3 3 2 2 2 2 3 1 3 3 3 2 2 2 2 2 2 3 3 2 2 3 2 1 1
## [4353] 3 3 2 1 2 2 3 3 3 2 3 3 3 2 3 2 3 3 1 2 3 2 2 3 3 2 2 3 3 2 2 3 3 2
## [4387] 3 1 3 3 2 1 2 2 2 3 3 3 3 2 1 1 3 1 3 3 3 2 3 3 1 3 3 1 2 3 3 2 1 1 1
## [4421] 3 2 3 2 3 1 3 2 2 3 2 2 2 2 3 2 2 3 2 2 1 3 3 2 3 3 3 2 1 3 2 1 3 2
## [4455] 3 3 1 2 3 2 3 2 3 3 2 3 3 3 3 3 3 3 3 2 3 3 3 2 2 2 1 3 1 2 3 3 3
## [4489] 3 1 2 1 3 2 3 3 3 2 2 3 2 2 2 1 1 1 1 2 2 2 3 2 2 3 3 2 2 2 2 2 2 2
## [4523] 2 3 2 3 2 3 2 3 2 3 2 2 3 2 2 2 3 2 2 3 1 2 3 2 3 2 3 2 2 1 2 2 3
## [4557] 2 3 3 2 1 2 3 3 3 2 3 3 3 3 3 3 3 2 1 2 3 2 3 2 2 2 2 2 2 2 2 2 1
## [4591] 3 2 3 2 2 3 1 3 3 2 3 2 3 1 2 2 2 3 3 3 2 3 3 2 3 2 3 3 3 3 3 3 1 3
## [4625] 2 3 2 3 2 3 3 3 1 1 2 3 2 2 1 1 3 3 3 3 3 2 2 3 3 2 1 3 3 2 2 2 3 1
## [4659] 3 2 2 2 3 2 3 1 2 3 2 3 2 1 3 2 1 3 2 2 1 3 1 2 2 3 3 3 3 3 1 3 2 1
## [4693] 1 3 3 3 3 2 1 2 3 3 3 3 2 2 2 2 2 3 3 2 2 2 2 3 3 3 2 2 3 2 2 3 2 1
## [4727] 2 1 2 3 2 2 2 3 2 1 2 2 3 3 1 2 2 1 3 3 2 3 3 2 3 1 2 3 2 2 2 3 1 3
## [4761] 3 2 3 1 2 2 3 3 3 3 3 3 3 3 1 2 3 3 2 2 2 3 2 3 3 2 2 2 1 2 2 3 2 2
## [4795] 2 2 2 2 2 3 2 2 2 2 2 2 2 2 3 3 3 2 3 2 2 3 1 3 2 2 2 3 3 3 1 3 2 3
## [4829] 3 2 2 2 3 2 2 3 1 2 2 3 2 2 1 3 2 3 2 1 2 2 2 3 3 3 2 3 3 1 1 3 3 2
## [4863] 3 1 3 2 3 2 1 1 3 3 3 1 2 3 3 2 2 3 3 2 3 2 3 3 1 3 2 2 2 2 3 2 2 3
## [4897] 3 3 2 2 3 2 3 3 2 3 3 3 1 2 3 1 3 3 2 2 2 1 2 2 2 2 2 2 3 3 3 3 2 1
## [4931] 3 1 3 2 3 2 2 2 2 2 1 2 3 2 3 1 3 3 2 1 1 3 2 1 1 2 3 2 3 2 3 2 3 2
## [4965] 2 1 1 2 2 3 3 2 3 3 3 1 2 3 3 3 2 2 3 2 3 3 1 2 3 3 2 3 2 3 1 2 3 2
## [4999] 2 2 2 3 3 3 2 3 3 2 2 1 2 1 3 2 2 3 2 3 1 2 3 3 2 2 2 1 3 3 3 2 2
## [5033] 2 2 2 2 3 1 3 3 3 3 3 1 3 2 2 2 3 2 1 1 3 2 1 1 3 2 3 2 2 3 2 1 2 2
## [5067] 1 2 1 2 3 2 3 2 3 2 3 2 3 3 2 2 3 3 1 2 1 2 2 3 2 3 1 3 3 2 3 2 1 2
## [5101] 2 2 3 3 2 3 2 2 3 3 3 2 2 2 2 3 3 2 2 1 3 2 2 3 1 2 2 2 2 2 3 2 3 2
## [5135] 2 1 3 2 2 3 1 3 3 1 2 2 2 2 3 3 2 2 2 1 1 2 2 3 2 2 3 3 3 2 3 1 2 3
## [5169] 1 2 3 3 2 2 1 2 2 3 3 3 2 2 1 2 2 3 2 3 3 1 3 3 3 2 3 3 2 3 1 3 3 3
## [5203] 3 3 3 2 2 3 3 3 3 2 3 2 2 3 3 2 2 3 2 3 2 2 2 2 2 3 2 1 3 3 2 3 2 2
## [5237] 2 2 3 2 2 2 2 3 2 1 3 2 2 2 3 3 2 2 2 3 3 2 3 3 3 3 2 2 3 3 1 2 3 2
## [5271] 1 2 2 3 2 2 3 3 3 2 2 1 1 2 2 3 3 3 2 2 1 3 3 3 3 2 3 2 1 2 2 1 2 2
## [5305] 1 1 3 1 2 1 1 3 2 2 2 2 3 2 1 3 2 3 3 3 3 2 1 3 3 2 3 3 3 3 2 1 3 2
## [5339] 2 2 2 3 2 3 2 3 3 3 2 3 2 3 2 3 1 3 3 2 3 3 1 2 1 1 2 1 2 3 2 3 2 2
## [5373] 3 2 3 1 3 2 2 2 3 2 3 3 3 2 2 2 3 3 2 2 3 1 3 1 1 3 3 3 2 2 2 2 2 2
## [5407] 3 3 3 2 3 2 3 3 2 3 2 2 3 3 2 3 3 2 3 2 2 2 1 3 1 3 2 3 3 3 3 3 1 3
## [5441] 3 2 3 1 3 2 1 3 2 3 3 1 3 1 3 2 3 2 2 2 3 3 3 2 3 2 2 3 3 3 1 3 1 1
## [5475] 1 3 3 3 3 3 2 2 3 3 2 2 3 3 1 3 2 3 3 3 2 2 3 2 3 3 2 3 2 3 1 3 3 2
## [5509] 2 3 2 3 3 3 2 3 2 2 2 3 2 3 2 3 3 2 2 3 3 1 2 2 2 1 3 3 3 3 3 3 3 2
## [5543] 3 3 3 2 3 3 1 2 2 2 2 2 2 3 3 2 2 2 3 3 2 3 3 2 3 2 2 2 3 1 3 3 2 2 3 2
## [5577] 2 3 2 2 2 3 3 3 2 2 2 3 2 3 3 3 3 2 3 3 3 2 2 3 3 2 3 3 3 1 3 3 3 3
## [5611] 1 3 2 1 2 1 2 2 3 2 2 2 2 3 2 3 2 1 2 3 3 3 3 3 2 2 2 1 2 3 3 2 2
## [5645] 3 2 2 3 3 3 2 2 3 3 3 3 2 3 2 3 3 3 3 2 2 3 3 2 3 2 3 2 2 3 3 3 3 2
## [5679] 3 2 2 3 3 3 2 2 3 3 2 3 3 3 1 2 2 1 2 2 2 2 2 2 3 2 3 2 2 2 3 2 3 3 3
## [5713] 2 3 2 2 3 1 1 2 2 2 2 1 1 2 2 3 3 3 2 3 2 2 1 2 2 3 2 2 1 3 3 2 1 2
## [5747] 2 3 3 1 3 3 3 3 3 3 1 3 2 2 2 2 3 2 3 3 2 3 2 3 3 2 3 1 3 3 2 3 3 3
## [5781] 2 3 3 2 3 3 2 3 3 2 3 3 3 3 3 2 3 2 3 2 3 3 3 2 3 1 2 3 2 3 2 1 3 3 2
## [5815] 1 1 2 3 3 2 2 3 2 3 3 2 3 2 3 3 3 2 3 3 3 3 1 1 2 2 3 2 2 3 3 2 3 3
## [5849] 1 3 1 1 2 2 3 3 2 3 3 2 3 3 2 2 3 1 3 3 2 3 3 3 1 3 2 3 3 3 2 3 1 2
## [5883] 1 2 3 2 3 3 2 1 3 3 2 3 3 3 3 2 2 2 3 3 1 2 3 2 3 1 3 3 2 3 2 1 1 3
## [5917] 3 3 2 3 2 2 3 2 3 2 2 1 3 2 3 3 3 3 2 2 1 3 2 3 1 1 3 3 2 3 2 3 3 3
```

```
## [5951] 1 2 2 3 3 2 3 2 2 2 3 2 3 3 3 3 3 3 3 2 3 3 3 3 3 3 2 3 2 2 2 2 2 3 1 2
## [5985] 2 3 2 2 2 2 2 2 2 2 3 2 1 2 3 2 3 2 1 3 1 3 2 1 2 3 3 3 3 2 3 1 3 3 2
## [6019] 3 2 2 3 2 2 3 2 3 2 3 3 3 2 3 3 3 3 1 2 2 3 2 1 2 3 2 2 3 2 2 3 2 3
## [6053] 3 2 3 2 2 1 3 3 3 2 3 3 1 2 2 2 3 2 1 3 2 3 3 2 3 3 1 3 3 3 2 3 1 3
## [6087] 3 2 2 2 2 2 1 3 2 1 3 2 2 2 3 1 3 2 2 2 2 3 2 3 3 3 1 1 3 3 2 3 2 2
## [6121] 3 3 3 3 1 2 2 3 1 3 3 2 3 2 3 1 2 2 1 2 3 2 2 3 3 2 2 1 3 2 1 3 3 2
## [6155] 1 2 1 3 3 2 2 3 3 2 2 2 2 3 3 2 2 3 2 2 3 2 2 2 3 2 3 3 3 3 3 2 3 2
## [6189] 2 2 2 3 2 3 3 2 3 3 1 3 2 3 3 1 2 3 3 3 3 2 3 1 3 2 3 3 3 3 2 3 3 3
## [6223] 3 3 2 2 2 2 3 3 3 2 3 3 3 3 2 1 3 2 3 2 2 2 2 3 2 3 3 3 2 3 2 2 3 3
## [6257] 2 2 2 2 2 3 2 3 2 2 2 2 2 3 2 3 2 2 2 3 2 2 2 3 2 2 2 3 2 1 2 3 1
## [6291] 1 1 3 2 3 3 2 2 2 3 2 2 3 2 3 3 3 3 3 3 3 3 1 3 2 2 2 3 3 2 3 2 3 2
## [6325] 2 2 2 3 2 3 3 3 3 2 2 2 2 3 3 3 2 3 2 3 2 2 3 2 3 1 2 2 2 1 1 1 3 2
## [6359] 3 3 3 2 3 3 3 2 3 2 2 2 3 3 1 3 2 2 2 3 2 3 3 2 2 3 1 2 3 3 3 3 2 3
## [6393] 2 3 2 2 3 2 3 3 3 2 3 1 2 3 2 2 3 1 2 2 3 3 2 3 3 2 3 3 3 3 2 3 3 1
## [6427] 2 3 1 3 3 3 2 1 3 3 3 2 2 3 2 2 2 1 2 1 3 2 2 2 2 2 3 3 1 3 3 2 3 1
## [6461] 2 3 3 3 1 3 3 2 3 3 3 3 1 2 2 2 2 3 2 3 2 1 3 2 3 3 3 1 3 2 3 1 3 2
## [6495] 2 3 3 3 2 2 2 2 3 2 2 1 3 1 2 1 2 2 3 2 3 3 2 2 2 2 3 3 2 2 3 1 3 3
## [6529] 2 2 3 2 2 3 2 3 2 3 1 2 2 2 2 3 2 3 2 2 3 3 3 3 2 3 1 2 2 2 3 2 2 3
## [6563] 1 2 3 1 2 2 2 2 2 2 2 2 2 2 3 2 2 3 2 3 3 2 2 3 2 3 2 2 2 3 3 2 1 1
## [6597] 1 3 2 1 2 1 2 3 2 3 3 2 3 3 3 2 3 2 1 2 3 2 3 3 3 3 2 3 2 2 3 2 2 2
## [6631] 3 1 2 2 3 3 3 3 3 2 2 3 3 2 3 3 2 2 3 2 2 3 1 3 2 3 3 3 2 3 2 2 2 1
## [6665] 2 3 3 3 2 2 3 3 2 1 3 3 2 2 3 3 2 2 3 3 3 3 3 3 1 2 2 3 2 3 3 2 2 2 3
## [6699] 3 3 3 2 3 3 2 2 3 3 2 3 3 2 2 2 3 3 3 3 2 3 3 3 2 1 2 3 3 2 3 1 3 2
## [6733] 2 3 3 2 2 3 2 2 2 2 2 1 2 3 1 1 2 2 2 2 3 2 2 3 3 3 3 2 2 3 1 3 2 1
## [6767] 2 1 2 2 3 3 3 3 1 3 2 2 2 2 2 3 3 2 3 2 2 1 2 1 3 3 3 2 2 3 3 3 3 2
## [6801] 3 2 2 2 3 2 3 2 2 3 2 3 3 2 3 3 3 2 2 2 3 2 2 3 1 2 3 2 3 3 2 1 3 3
## [6835] 3 3 3 2 1 3 2 3 3 3 3 3 3 3 3 3 1 1 3 3 1 3 3 2 2 3 3 2 2 2 3 3 2 3
## [6869] 2 3 2 3 2 3 3 2 3 3 3 3 3 3 3 1 2 2 1 2 2 3 3 3 2 3 2 2 1 2 3 2 2 2
## [6903] 2 2 3 1 1 2 2 2 2 3 3 1 1 3 3 2 3 3 3 2 1 1 1 2 3 2 3 3 1 2 1 2 2 2
## [6937] 2 3 2 2 2 3 2 3 3 2 2 3 3 3 2 3 3 3 1 3 3 3 3 2 1 3 2 2 2 3 3 3 1 3
## [6971] 3 2 3 2 3 3 2 2 1 1 3 3 2 1 3 2 3 1 1 2 2 2 2 2 3 3 3 2 3 3 3 3 2 1
## [7005] 2 1 2 2 3 1 3 2 3 1 2 2 3 1 3 2 2 3 2 2 3 2 3 3 3 3 1 3 2 3 3 2 3 2
## [7039] 3 3 1 2 2 3 3 2 3 2 2 2 2 2 1 3 1 3 2 3 3 3 3 3 2 3 2 1 2 2 1 2 2 3 3
## [7073] 2 2 2 2 2 2 3 1 3 3 3 2 3 2 3 2 2 3 2 3 2 2 2 3 2 3 2 2 2 2 3 2 3 2
## [7107] 2 2 2 2 2 2 3 1 3 2 1 3 1 2 2 2 3 2 3 1 2 3 2 2 1 2 2 2 2 3 3 2 2 2
## [7141] 3 3 2 2 2 1 3 1 3 3 3 2 2 2 3 3 2 2 2 3 2 3 1 3 2 2 2 3 2 2 1 3 2 2
## [7175] 3 2 2 2 3 2 3 3 2 3 3 3 3 3 3 3 1 2 3 3 1 2 2 2 3 2 3 1 3 1 1 1 2 3
## [7209] 2 3 3 2 2 1 2 3 2 2 3 2 2 2 3 3 2 1 3 2 3 2 2 3 2 2 1 3 3 3 3 2 2 3
## [7243] 3
##
## Within cluster sum of squares by cluster:
## [1]  3.169202  7.115763 16.140443
##  (between_SS / total_SS =  70.5 %)
##
## Available components:
##
## [1] "cluster"     "centers"     "totss"       "withinss"
## [5] "tot.withinss" "betweenss"   "size"        "iter"
## [9] "ifault"
```

#Step 5: Visualize the cluster data based on analysis above.

```r
# Vlisualize the cluster that I have performed above with fviz_cluster function.
fviz_cluster(final, data = df, geom="point")
```

**Cluster plot**



```r
#Convert data to geospatial format with sf library
crimes_sf <- st_as_sf(NYPD_shooting_incident,
                      coords = c("Longitude", "Latitude"),
                      crs = 4326)

#Plot geospatial location (lat,lon) and project to map with ggplot and geom_sf
class(crimes_sf)
```

```
## [1] "sf"          "tbl_df"      "tbl"          "data.frame"
```

```r
ggplot() +
  geom_sf(data = crimes_sf)
```

#Step 6: Identify Bias. Please see the analysis below. I have identified the bias on the data set: - Data set highly bias to race. Black people are very high compare with others in PERP_RACE. - Also the same with PERP_SEX, the perpetrator mostly are male.

```
#Check bias on the data
NYPD_shooting_incident %>% count(PERP_RACE, sort = TRUE)
```

```
## # A tibble: 7 x 2
##    PERP_RACE                          n
##    <chr>                          <int>
## 1 BLACK                           4938
## 2 WHITE HISPANIC                   861
## 3 UNKNOWN                          739
## 4 BLACK HISPANIC                   493
## 5 WHITE                            147
## 6 ASIAN / PACIFIC ISLANDER          64
## 7 AMERICAN INDIAN/ALASKAN NATIVE     1
```

```
NYPD_shooting_incident %>% count(PERP_SEX, sort = TRUE)
```

```
## # A tibble: 3 x 2
##    PERP_SEX      n
##    <chr>     <int>
## 1 M          6458
## 2 U           591
## 3 F           194
```

```
#Visualize it
ggplot(NYPD_shooting_incident, aes(x=as.factor(PERP_RACE), fill=as.factor(PERP_RACE) )) +
  geom_bar( ) +
  scale_fill_brewer(palette = "Set1") +
  theme(legend.position="none")
```



```
#Visualize it
ggplot(NYPD_shooting_incident, aes(x=as.factor(PERP_SEX), fill=as.factor(PERP_SEX) )) +
  geom_bar( ) +
  scale_fill_brewer(palette = "Set1") +
  theme(legend.position="none")
```