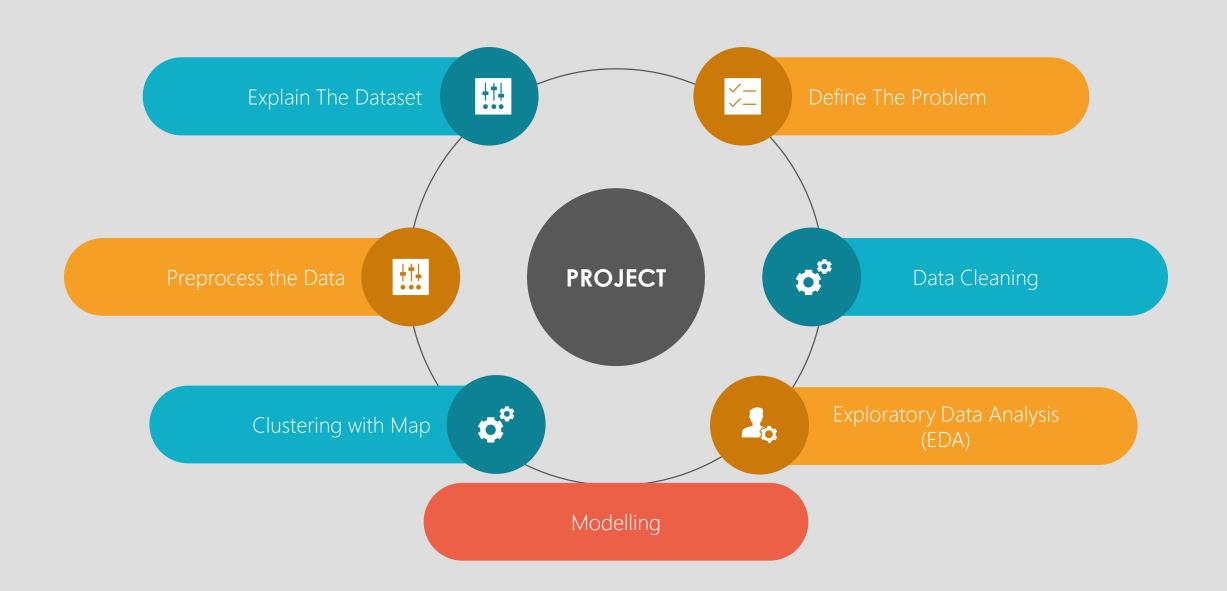
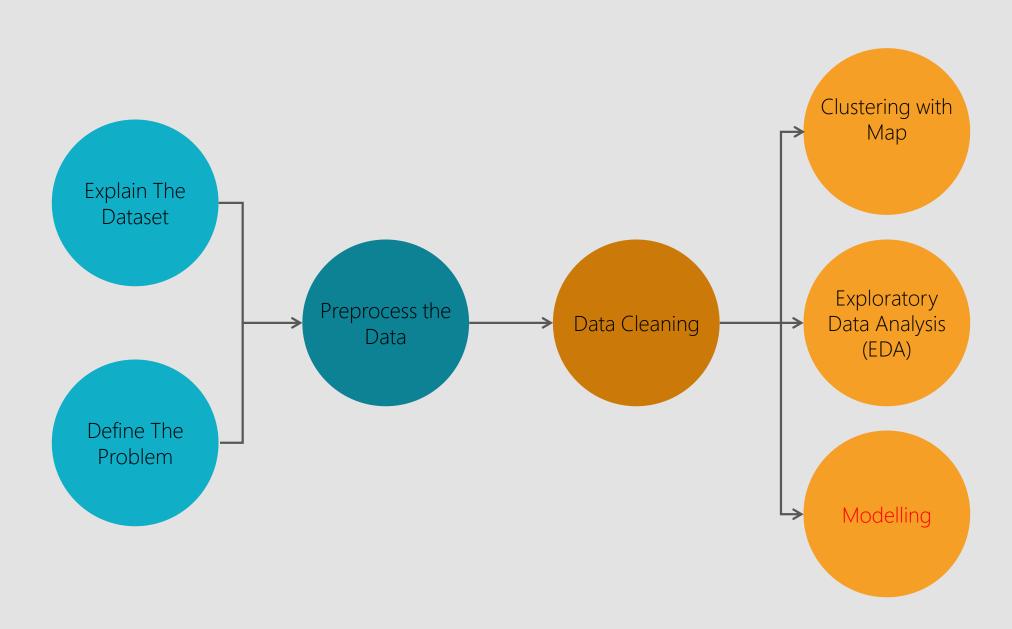


NY City Vehicle Analysis Statistical Learning Presentation

Project Analysis



Project Analysis



Explain The Dataset •



The Motor Vehicle Collisions crash table contains details on the crash event.

The Motor Vehicle Collisions data tables contain information from all police-reported motor vehicle collisions in NYC. The police report (MV104-AN) is required to be filled out for collisions where someone is injured or killed, or where there is at least \$1000 worth of damage

A brief summary of data with head() function

```
2019-08-05T00:00:00.000
                                16:30
                                                   11434 40.67605 -73.79018 {'type': 'Point', 'coordinates': [-73.790184, 40.676052]}
                                         QUEENS
2019-08-27T00:00:00.000
                                                   11225 40.65778 -73.95110 {'type': 'Point', 'coordinates': [-73.951096, 40.65778]}
                                       BROOKLYN
2019-08-15T00:00:00.000
                                17:57 MANHATTAN
                                                   10002 40.71814 -73.99384 {'type': 'Point', 'coordinates': [-73.993835, 40.718143]} CHRYSTIE STREET
                                                                                                                                                                               GRAND STREET
2019-08-30T00:00:00.000
                                           BRONX
                                                   10460 40.84053 -73.86661 {'type': 'Point', 'coordinates': [-73.86661, 40.840534]}
2019-08-06T00:00:00.000
                                 9:45 MANHATTAN
                                                                                {'type': 'Point', 'coordinates': [-73.9754, 40.74544]} EAST 35 STREET
                                                                                                                                                                                   2 AVENUE
                                                    11222 40.72652 -73.94639 {'type': 'Point', 'coordinates': [-73.94639, 40.726524]}
2019-08-29T00:00:00.000
150-08
          123 AVENUE
          HAWTHORNE STREET
1837
          EAST TREMONT AVENUE
          JEWEL STREET
                                                                                                                                  Unspecified
                                                                                          Passing Too Closely
                                                                                          Passing Too Closely
                                                                                                                                  Unspecified
                                                                               Driver Inattention/Distraction
                                                                                                   Unspecified
                                                                                                                                  Unspecified
                                                                             0 Driver Inattention/Distraction Driver Inattention/Distraction
                                                                                                  Unspecified
                                                                                                                                  Unspecified
                                                                                          VEHICLE.TYPE.CODE.1
                                                                                                                               VEHICLE.TYPE.CODE.2 VEHICLE.TYPE.CODE.3 VEHICLE.TYPE.CODE.4
                                                                                                                                     Pick-up Truck
                                                                  4195773 Station Wagon/Sport Utility Vehicle Station Wagon/Sport Utility Vehicle
                                                                                                         Taxi Station Wagon/Sport Utility Vehicle
                                                                  4198749
                                                                  4183798 Station Wagon/Sport Utility Vehicle
                                                                  4196772 Station Wagon/Sport Utility Vehicle
                                                                                                                                              Bike
VEHICLE.TYPE.CODE.5
```

Objectives

/_

- 1. The Problem: Big amount of accidents on daily basis
- 2. Download data on the kaggle
- 3. Import the data to R Studio
- 4. Visualize accidents and locations.
- 5. Select the variables in the dataset from the dataset as required. (date/location/time etc...)
- 6. Solution: Estimate the potential traffic accident area based on current locations.

With data <- data_try %>% select(ACCIDENT.DATE, ACCIDENT.TIME, LATITUDE, LONGITUDE)

We select our variables

All 40.676032 40.676032 40.676032 40.676032 40.6778 40.718143 40.840534 40.726524 40.667522 40.86821 40.686732	-73.9754 -73.94639 -73.78063	All 2019-08-05 2019-08-27 2019-08-15 2019-08-30 2019-08-29 2019-08-31 2019-08-11	All 2019-08-05T16-30:00Z 2019-08-25T16-30:00Z 2019-08-25T16-02:00Z 2019-08-15T17:57:00Z 2019-08-30T21:53:00Z 2019-08-06T09:45:00Z 2019-08-29T12:28:00Z 2019-08-31T02:16:00Z 2019-08-11T22:23:00Z 2019-08-11T22:	Meekday Pazartesi Sali Perşembe Cuma Sali Perşembe Cumatesi Pazar	Weekend All Weekdays Weekdays Weekdays Weekdays Weekdays Weekdays Weekdays Weekdays Weekdays	ф [АШ	17 17 18 22 10 13 3
40.67632 40.65778 40.718143 40.840534 40.74544 40.726524 40.667522 40.85821 40.666492	-73.790184 -73.951096 -73.993835 -73.86661 -73.9754 -73.94639 -73.78063 -73.78063	2019-08-05 2019-08-27 2019-08-15 2019-08-30 2019-08-30 2019-08-29 2019-08-31 2019-08-11	2019-08-05T16:30:00Z 2019-08-27T16:02:00Z 2019-08-15T17:57:00Z 2019-08-30T21:53:00Z 2019-08-06T09:45:00Z 2019-08-29T12:28:00Z 2019-08-31T02:16:00Z	Pazartesi Sali Perşembe Cuma Sali Perşembe Cumartesi	Weekdays Weekdays Weekdays Weekdays Weekdays Weekdays Weekdays	All	17 18 22 10
40.65778 40.718143 40.840534 40.74544 40.726524 40.667522 40.85821 40.666492	-73.951096 -73.993835 -73.86661 -73.9754 -73.94639 -73.78063 -73.91679	2019-08-27 2019-08-15 2019-08-30 2019-08-06 2019-08-29 2019-08-31 2019-08-11	2019-08-27T16-02:00Z 2019-08-15T17:57:00Z 2019-08-30T21:53:00Z 2019-08-06T09:45:00Z 2019-08-29T12:28:00Z 2019-08-31T02:16:00Z	Salı Perşembe Cuma Salı Perşembe Cumartesi	Weekdays Weekdays Weekdays Weekdays Weekdays		17 18 22 10
40.718143 40.840534 40.74544 40.726524 40.667522 40.85821 40.666492	-73,993835 -73,86661 -73,9754 -73,94639 -73,78063 -73,91679	2019-08-15 2019-08-30 2019-08-06 2019-08-29 2019-08-31 2019-08-11	2019-08-30T21:53:00Z 2019-08-30T21:53:00Z 2019-08-06T09:45:00Z 2019-08-29T12:28:00Z 2019-08-31T02:16:00Z	Perşembe Cuma Salı Perşembe Cumartesi	Weekdays Weekdays Weekdays Weekdays Weekdays		18 22 10 13
40.840534 40.74544 40.726524 40.667522 40.85821 40.666492	-73.86661 -73.9754 -73.94639 -73.78063 -73.91679	2019-08-30 2019-08-06 2019-08-29 2019-08-31 2019-08-11	2019-08-30T21:53:00Z 2019-08-06T09:45:00Z 2019-08-29T12:28:00Z 2019-08-31T02:16:00Z	Cuma Salı Perşembe Cumartesi	Weekdays Weekdays Weekdays		22 10 13
40.74544 40.726524 40.667522 40.85821 40.666492	-73.9754 -73.94639 -73.78063 -73.91679	2019-08-06 2019-08-29 2019-08-31 2019-08-11	2019-08-06T09:45:00Z 2019-08-29T12:28:00Z 2019-08-31T02:16:00Z	Salı Perşembe Cumartesi	Weekdays Weekdays Weekdays		10 13
40.726524 40.667522 40.85821 40.666492	-73.94639 -73.78063 -73.91679	2019-08-29 2019-08-31 2019-08-11	2019-08-29T12:28:00Z 2019-08-31T02:16:00Z	Perşembe Cumartesi	Weekdays Weekdays		13
40.667522 40.85821 40.666492	-73.78063 -73.91679	2019-08-31 2019-08-11	2019-08-31T02:16:00Z	Cumartesi	Weekdays		
40.85821 40.666492	-73.91679	2019-08-11					3
40.666492			2019-08-11T22:23:00Z	Pazar	Weekend		
	-73.76536						23
40.02772		2019-08-02	2019-08-02T16:30:00Z	Cuma	Weekdays		17
40.83772	-73.92763	2019-08-21	2019-08-21T17:30:00Z	Çarşamba	Weekdays		18
40.68164	-73.98568	2019-08-20	2019-08-20T14:25:00Z	Salı	Weekdays		15
40.644047	-73.95747	2019-08-17	2019-08-17T14:02:00Z	Cumartesi	Weekdays		15
40.698463	-73.960205	2019-08-20	2019-08-20T10:30:00Z	Salı	Weekdays		11
40.716652	-73.8259	2019-08-04	2019-08-04T23:09:00Z	Pazar	Weekend		24
40.67855	-73.81422	2019-08-12	2019-08-12T11:20:00Z	Pazartesi	Weekdays		12
	40.644047 40.698463 40.716652	40.644047 -73.95747 40.698463 -73.960205 40.716652 -73.8259	40.644047 -73.95747 2019-08-17 40.698463 -73.960205 2019-08-20 40.716652 -73.8259 2019-08-04	40.644047 -73.95747 2019-08-17 2019-08-17T14:02:00Z 40.698463 -73.960205 2019-08-20 2019-08-20T10:30:00Z 40.716652 -73.8259 2019-08-04 2019-08-04T23:09:00Z	40.644047 -73.95747 2019-08-17 2019-08-17T11-02:00Z Cumartesi 40.698463 -73.960205 2019-08-20 2019-08-20T10:30:00Z Salı 40.716652 -73.8259 2019-08-04 2019-08-04T23-09:00Z Pazar 40.67855 -73.81422 2019-08-12 2019-08-12T11:20:00Z Pazartesi	40.644047 -73.95747 2019-08-17 2019-08-17114-02-00Z Cumartesi Weekdays 40.698463 -73.960205 2019-08-20 2019-08-20710-30-00Z Salı Weekdays 40.716652 -73.8259 2019-08-04 2019-08-04T23:09-00Z Pazar Weekend 40.67855 -73.81422 2019-08-12 2019-08-12T11:20-00Z Pazartesi Weekdays	40.6464047 -73.95747 2019-08-17 2019-08-17T14-02-00Z Cumartesi Weekdays 40.698463 -73.960205 2019-08-20 2019-08-20T10-30-00Z Salı Weekdays 40.716652 -73.8259 2019-08-04 2019-08-04T23-09-00Z Pazar Weekend 40.67855 -73.81422 2019-08-12 2019-08-12T11-20-00Z Pazartesi Weekdays

widtharpoonup Summary of The Data widtharpoonup



```
> data %>% summary()
 ACCIDENT.DATE
                   ACCIDENT.TIME
                                         LATITUDE
                                                         LONGITUDE
 Length: 1612178
                   Length: 1612178
                                      Min.
                                             : 0.00
                                                       Min.
                                                              :-201.24
 Class : character
                   Class : character
                                      1st Qu.:40.67
                                                       1st Qu.: -73.98
                                      Median :40.72
                                                       Median : -73.93
       :character
                   Mode :character
Mode
                                             :40.69
                                                              : -73.87
                                      Mean
                                                       Mean
                                      3rd Qu.:40.77
                                                       3rd Qu.: -73.87
                                      Max.
                                             :42.32
                                                       Max.
                                                                  0.00
                                      NA's
                                             :196285
                                                       NA's
                                                              :196285
```

Data Cleaning



In fact, we wanted to do a data manipulation here because the location variables in the data were too large. For this we would use knn, random forest to do an undefined value tuning.

However, this process takes too long (because we do not have enough processing power). Here is a simple data editing

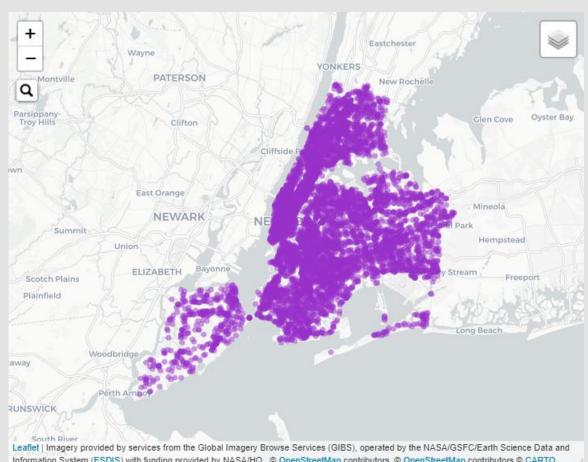
```
library (missForest)
                      summary (na data)
anyNA(cln na data)
                      rf source <- na data %>%
                        select (LATITUDE, LONGITUDE, NUMBER.OF. PERSONS.INJURED, NUMBER.OF. PERSONS.KILLED)
                 media rf data <- missForest(rf_source, ntree = 5)</pre>
                        <- sapply(na_data, function(x) which(is.na(x)))
                                                                                                             n)
                      212 1/2 data ((1512/1717)) FISTITUDENGITUDE < -72, LONGITUDE > -75)
anyNA(na data)
                                                                                                             ian)
                media cln na data[c(1$LONGITUDE),]$LONGITUDE
                      rf_data <- rf_data$ximp
                      rf data[c(1$LATITUDE),]$LATITUDE
mean data$NUMBER.OF.p|rf_data[c(1$LONGITUDE),]$LONGITUDE
mean data$NUMBER.OF.P
                                                                                                             OF. PERSONS.KILLED)
mean data$LATITUDE[is mean(cln_na_data[c(l$LATITUDE),]$LATITUDE - rf data[c(l$LATITUDE),]$LATITUDE)
mean data$LONGITUDE[i:summary(cln na data$LATITUDE)
                      mean(cln na data[c(1$LONGITUDE),]$LONGITUDE - rf data[c(1$LONGITUDE),]$LONGITUDE)
summary(mean data)
                      summary(cln na data$LONGITUDE)
```

Create a Map



When I look at the summary results,

We took the map on the right and it with 7500 samples.



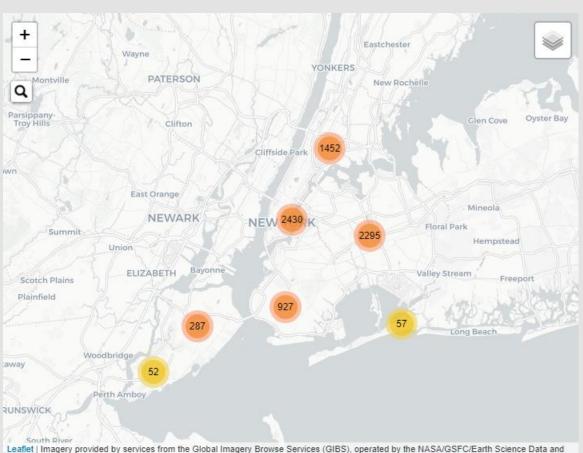
Information System (ESDIS) with funding provided by NASA/HQ., @ OpenStreetMap contributors, @ OpenStreetMap contributors

Map with Clustering



Thanks to that piece of code

We can organize our maps data.



Leaflet | Imagery provided by services from the Global Imagery Browse Services (GIBS), operated by the NASA/GSFC/Earth Science Data and Information System (ESDIS) with funding provided by NASA/HQ., @ OpenStreetMap contributors, @ OpenStreetMap contributors @ CARTO

Making and Converting



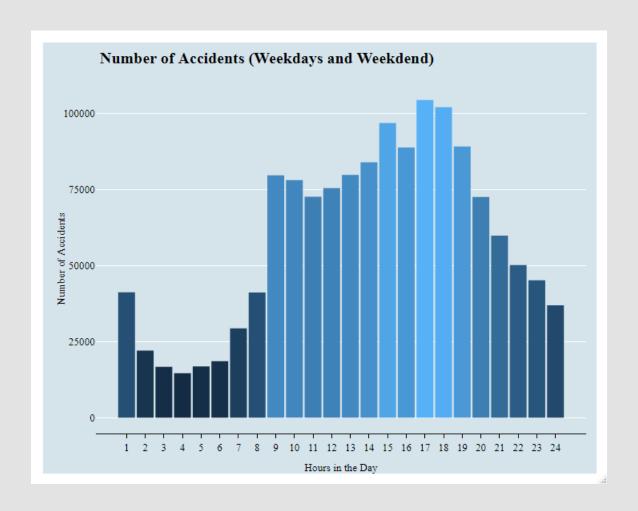
- We converted date time to time series data and created variables such as hour, weekday, weekend, etc.
- Hour is denoting 24 hours a day.
- However, we had to write the days in Turkish because the system language of the computer we applied was Turkish.

<u>So</u>: Saturday = Cumartesi, Sunday = Pazar,

Exploratory Data Analysis (EDA)



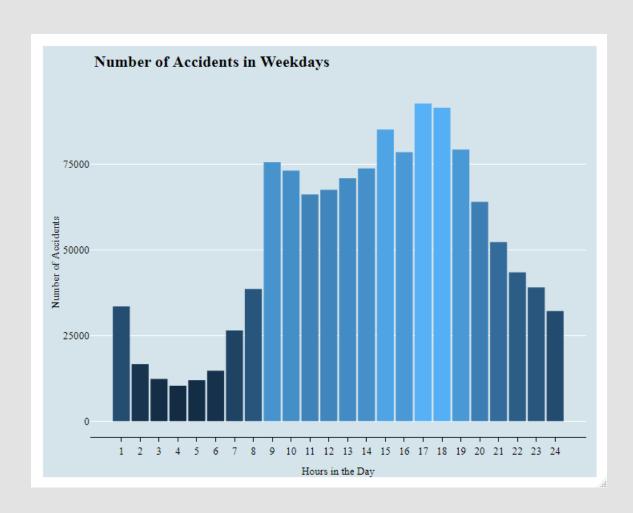
- When we look at the number of accidents according to the time schedule of the day, there is an increase in accidents almost regularly.
- However, we can see that this increase ended at 17 ~ 18 hours.



Accident in Weekdays



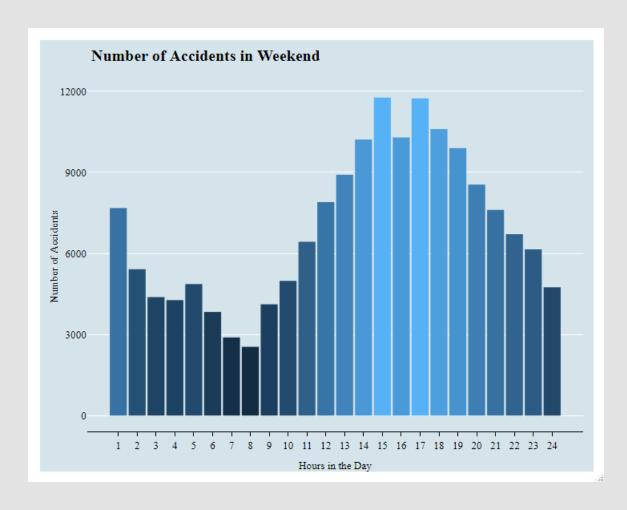
• It is the same as the observations from the data set, which is slightly higher in the morning.



Accident in Weekends



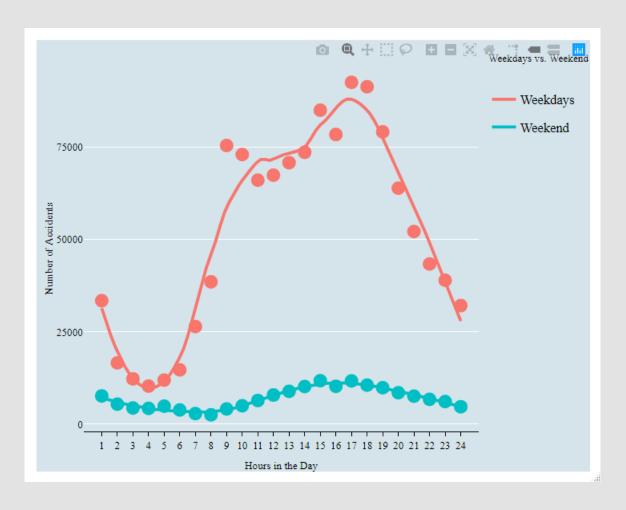
• For the weekend, the graph has changed in appearance, and the peak takes place between 15 and 17 hours.



Weekdays vs Weekend



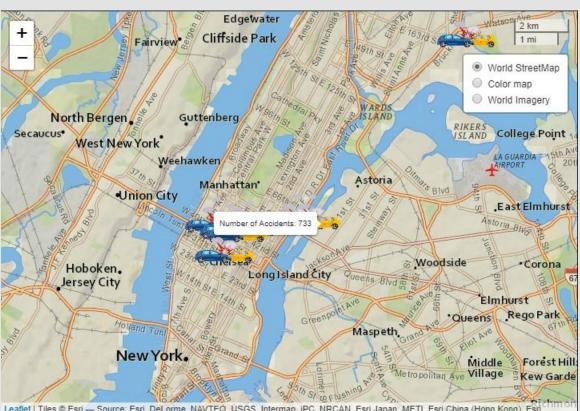
• If we compare the weekdays and weekend tables, we can say that people drive less risky vehicles on weekends.



Top 10 Dangerous Areas in Weekdays

10

- Instead of directing the calculation of the <u>top</u> <u>10 accident locations</u>, we processed the data a bit.
- If We use direct longitude and latitude data, the same collection point with slightly different coordinates is considered to be different collection locations and will definitely deviate from the actual result.
- Therefore, I rounded up to 3 decimal places without making a huge change to longitude and latitude. We also uploaded a picture to Ogulcan's blog to indicate the accident sites and used that picture as an icon. The graphic is interactive and can zoom in and out. If you place the mouse over the crashed car icon, it shows how many accidents are in the position according to the data set.



Leaflet | Tiles © Esri — Source: Esri, DeLorme, NAVTEQ, USGS, Intermap, iPC, NRCAN, Esri Japan, METI, Esri China (Hong Kong), Esri (Thailand), TomTom, 2012, Imagery provided by services from the Global Imagery Browse Services (GIBS), operated by the NASA/GSFC/Earth Science Data and Information System (ESDIS) with funding provided by NASA/HQ., © OpenStreetMap contributors, © CARTO, Tiles © Esri — National Geographic, Esri, DeLorme, NAVTEQ UNEP-WCMC, USGS, NASA, ESA, METI, NRCAN, GEBCO, NOAA, iPC

Top 10 Dangerous Areas in Weekends



• The locations of the vehicles, namely the regions where the most accidents occur, change directly on weekends and weekdays.



contributors @ CARTO, Tiles @ Esri — National Geographic, Esri, DeLorme, NAVTEQ/UNEP-WGMO, USGS, NASA, ESA, METI, NRCAN, GEBCO, NOAA, iPC



Leaflet | Tiles @ Esri - Source: Esri, DeLorme, NAVTEQ, USGS, Intermap, iPC, NRCAN, Esri Japan, METI, Esri China (Hong Kong), Esri (Thailand), TomTom, 2012, Imagery provided by services from the Global Imagery Browse Services (GIBS), operated by the NASA/GSFC/Earth Science Data and Information System (ESDIS) with funding provided by NASA/HQ. @ OpenStreetMap contributors. @ OpenStreetMap contributors @ CARTO, Tiles @ Esri - National Geographic, Esri, DeLorme, NAVTEQ, UNEP-WCMC, USGS, NASA, ESA, METI, NRCAN, GEBCO, NOAA, iPC

Modelling

Modelling

After examining all the accidents, we prepared a simple kmeans model to divide the accident points into 50 groups to suggest a potential accident point. With unsupervised learning, we got a simple but functional result.

```
data_coordinates <- data %>% select(LONGITUDE, LATITUDE)
data_fm <- data

'``{r}

set.seed(0)
data_kmeans <- data_coordinates %>%
    kmeans(50, nstart=20)
    save(data_kmeans, file = "input/kmeans_data.rda")
load("input/kmeans_data.rda")

data_fm%cluster <- data_kmeans$cluster

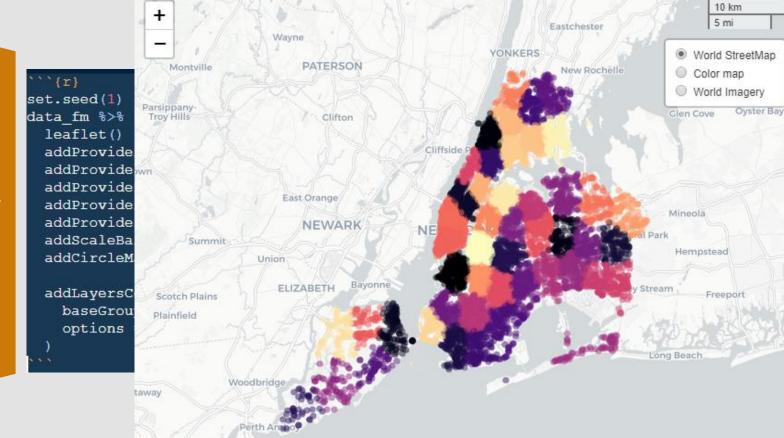
colors_in <- colorNumeric(
    palette = "magma",
    domain = data$cluster)</pre>
```

Modelling

RUNSWICK



Finally, we created a new car accident map using the kmeans model.



Leaflet | Tiles © Esri — Source: Esri, DeLorme, NAVTEQ, USGS, Intermap, iPC, NRCAN, Esri Japan, METI, Esri China (Hong Kong), Esri (Thailand), TomTom, 2012, Imagery provided by services from the Global Imagery Browse Services (GIBS), operated by the NASA/GSFC/Earth Science Data and Information System (ESDIS) with funding provided by NASA/HQ., © OpenStreetMap contributors, © OpenStreetMap contributors © CARTO

ry") %>%

New Top 10 **Estimated** Dangerous Area Map Weekdays

Modelling



New Top 10 Estimated Dangerous Area Map

Modelling

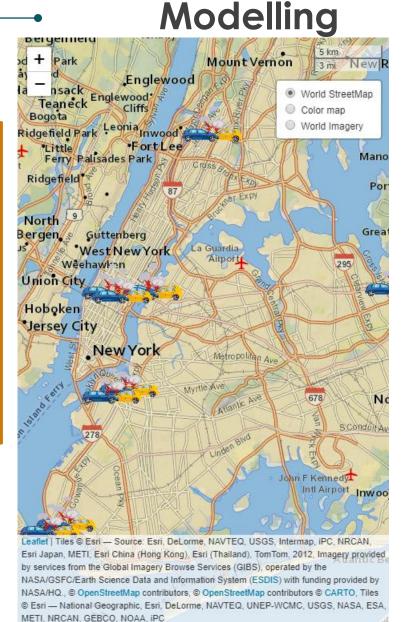


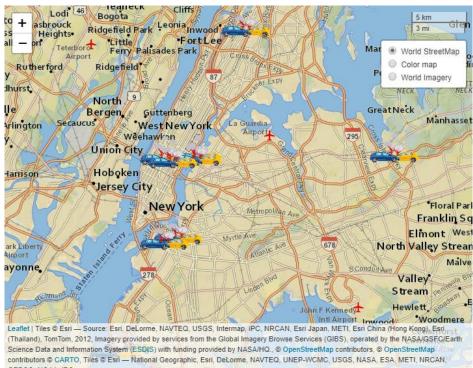


Esri Japan, METI, Esri China (Hong Kong), Esri (Thailand), TomTon, 2012, Imagery provided by services from the Global Imagery Browse Services (GIBS), operated by the NASA/GSFC/Earth Science Data and Information System (ESDIS) with funding provided by NASA/HQ., © OpenStreetMap contributors, © OpenStreetMap contributors, © CARTO, Tiles © Esri — National Geographic, Esri, DeLorme, NAVTEQ, UNEP-WCMC, USGS, NASA, ESA, METI, NRCAN, GEBCO, NOAA, iPC

METI, NRCAN, GEBCO, NOAA, iPC

Estimated Dangerous Area Map on Weekends vs Dangerous Area from data





References

- 1. https://www.kaggle.com/new-york-city/nypd-motor-vehicle-collisions
- 2. https://stackoverflow.com/questions/21382681/kmeans-quick-transfer-stage-steps-exceeded-maximum
- 3. https://github.com/tidyverse/lubridate/issues/669
- 4. http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf
- 5. https://stackoverflow.com/questions/49951416/how-to-use-colornumeric-within-addcircles-in-leaflet
- 6. https://medium.com/analytics-vidhya/a-guide-to-machine-learning-in-r-for-beginners-decision-trees-c24dfd490abb
- 7. https://medium.com/@mueller.johannes.j/use-r-and-gganimate-to-make-an-animated-map-of-european-students-and-their-year-abroad-517ad75dca06

