

# BALTIMORE Crimes



## ● Overview

**Problem:** Excessive occurrence of crime reported in Baltimore City

**Management:** Analyze the crime data via random forest and decision trees in order to find correlations in the data.

**Results:** Crime possibly correlates to different measure of time (day of week, month, etc.) and location.

**Next Steps:** Further analysis with better accuracy. Once completed: Implementation of guided tasks force; development of new protocols with respect to findings in time-related crime.

## ● Question

Can crime in Baltimore be correlated to certain measures of time and/or location?

If crime in Baltimore can be correlated to certain measures of time and/or location, then preventive methods could be installed to reduce crime based on this data

- **Presentation of the Dataset**

**276,202 crimes**

Recorded from Jan-2014 to Sept-2019

**Date of crime**

**Categories of crime**

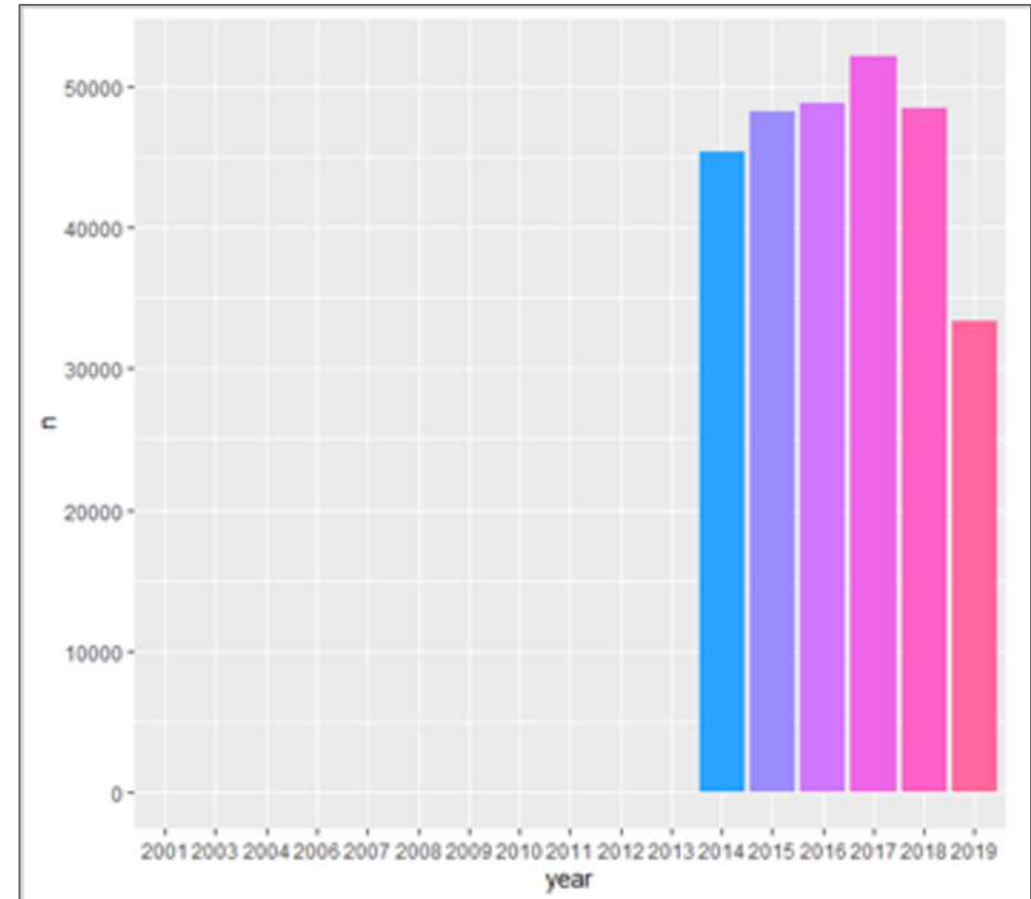
**Geographical area**

**Time of day**

# ● Sanitizing of the data

## Exploring the data

1. Inconsistent data before 2014
1. Information for Year 2019 only until 21st of September
1. Consolidate categories of crime (from 14 to 7)



## ● Data Exploration

Types of analysis/Analysis Methods:

- Plotting with ggplots
- Decision Tree
- Classification
  
- Alternative methods of analysis:
- Random forest

# ● Data Analysis

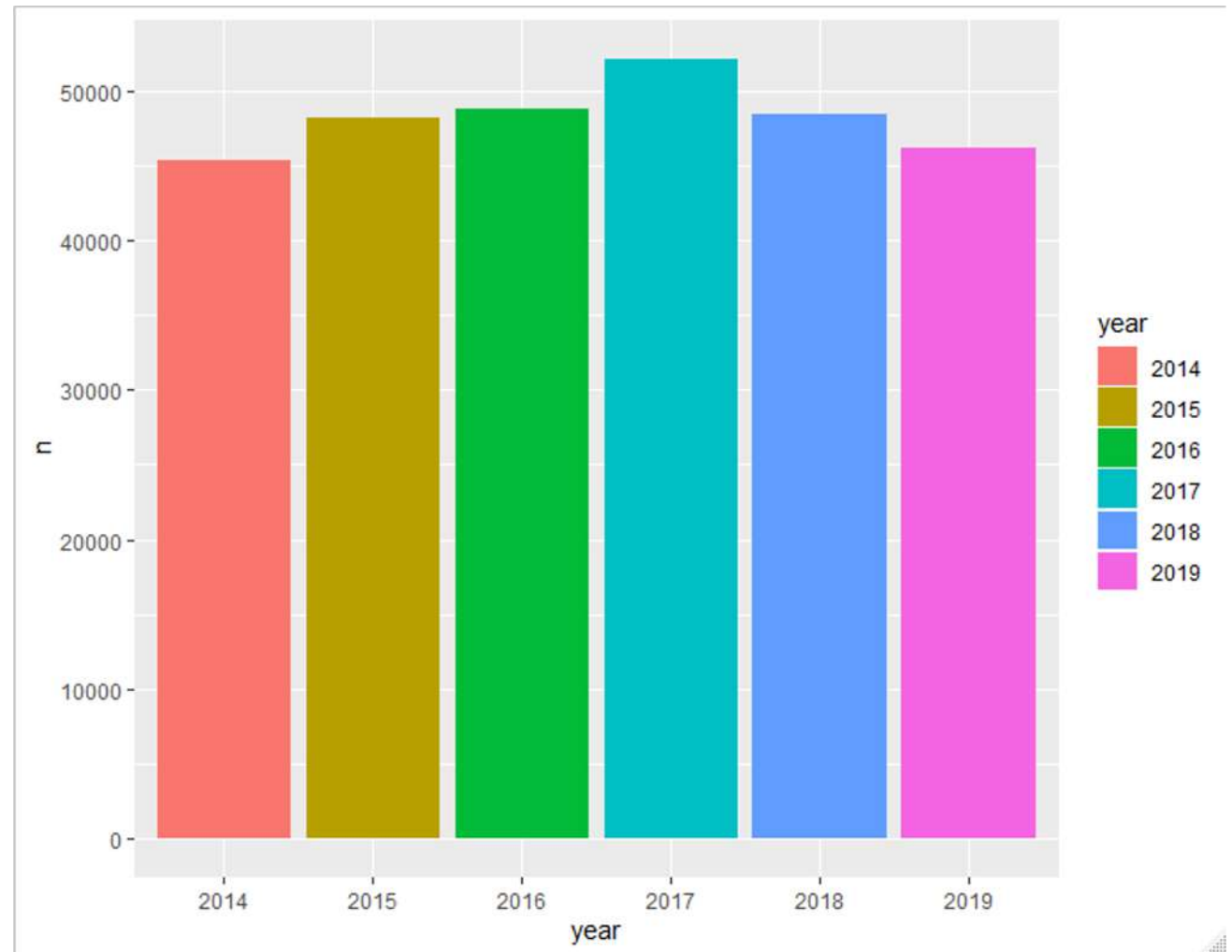
## Evolution of crimes through years

### Hypothesis:

- Correlation between years and total crimes

### Results:

- Total number of crimes of have at a steady pace, although not uniformly
- No immediate observable trend across the identified years



# ● Data Analysis

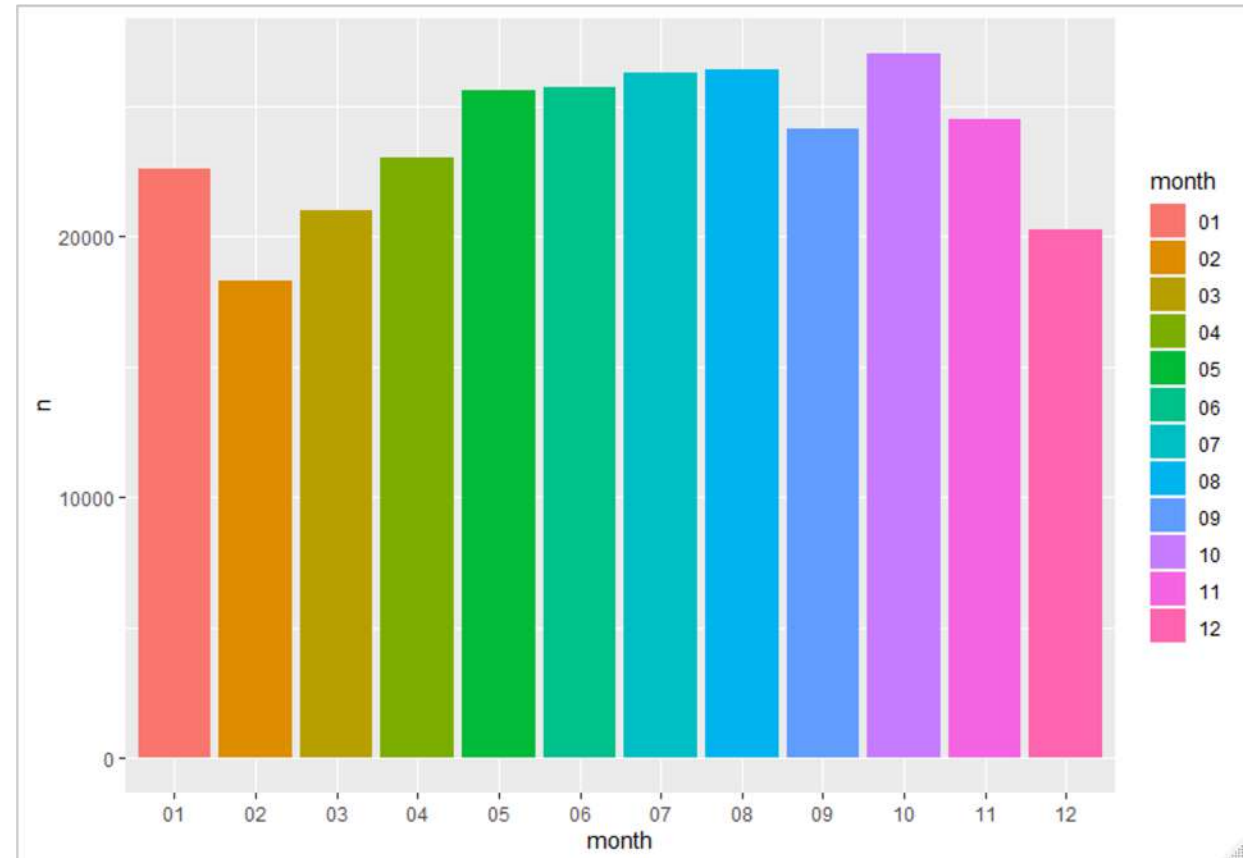
Evolution of crimes throughout the year

## Hypothesis:

- Correlation between the month and number of crimes

## Results:

- Months with most crimes: July and August
- Less crimes during winter than summer





# ● Data Analysis

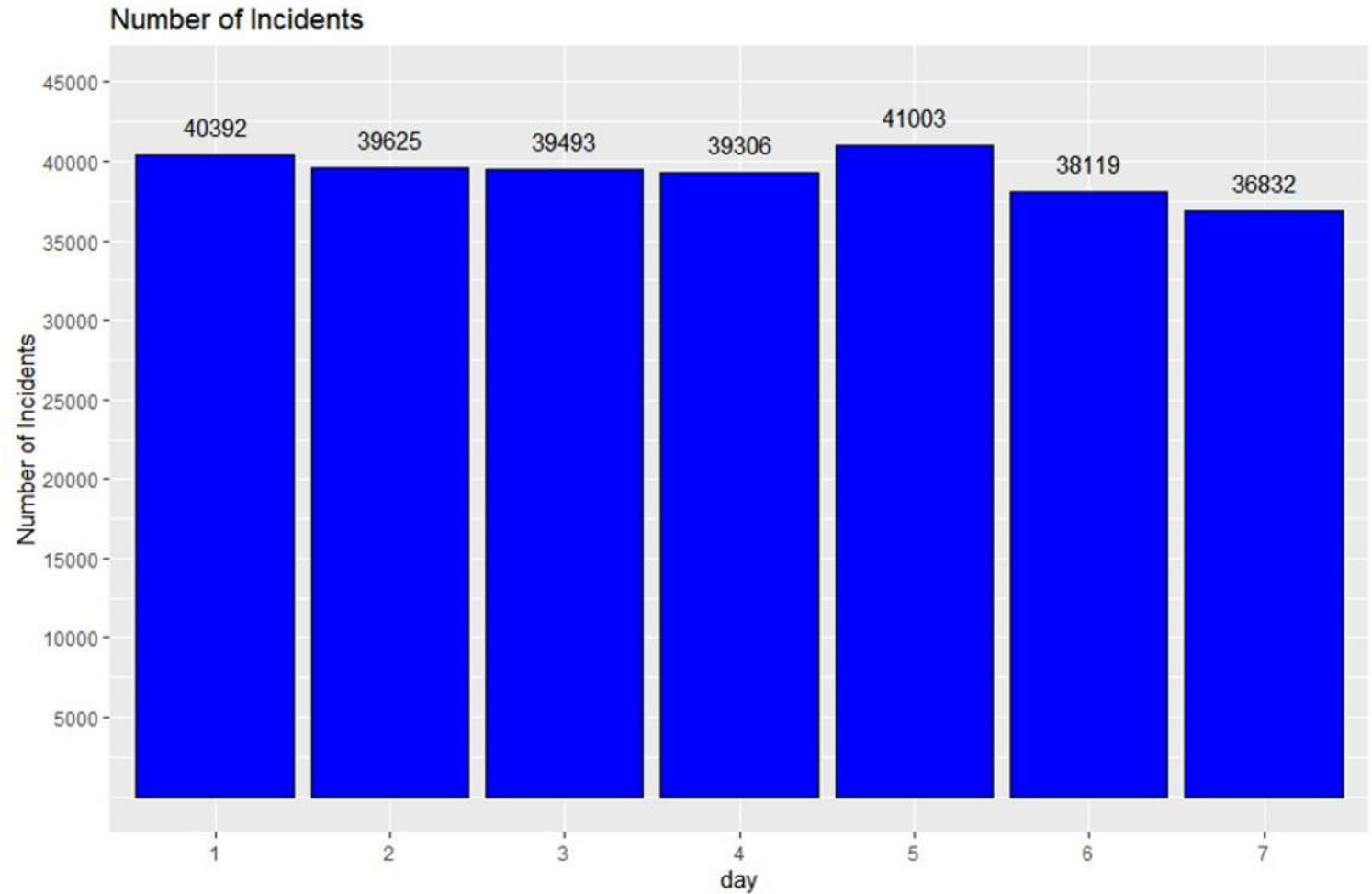
## Evolution of crimes throughout the week

### Hypothesis:

- Correlation between day of the week and number of crimes

### Results:

- Number of crimes almost constant
- Day with the most crimes: Friday
- Less crimes on the week-end than during the week



# ● Data Analysis

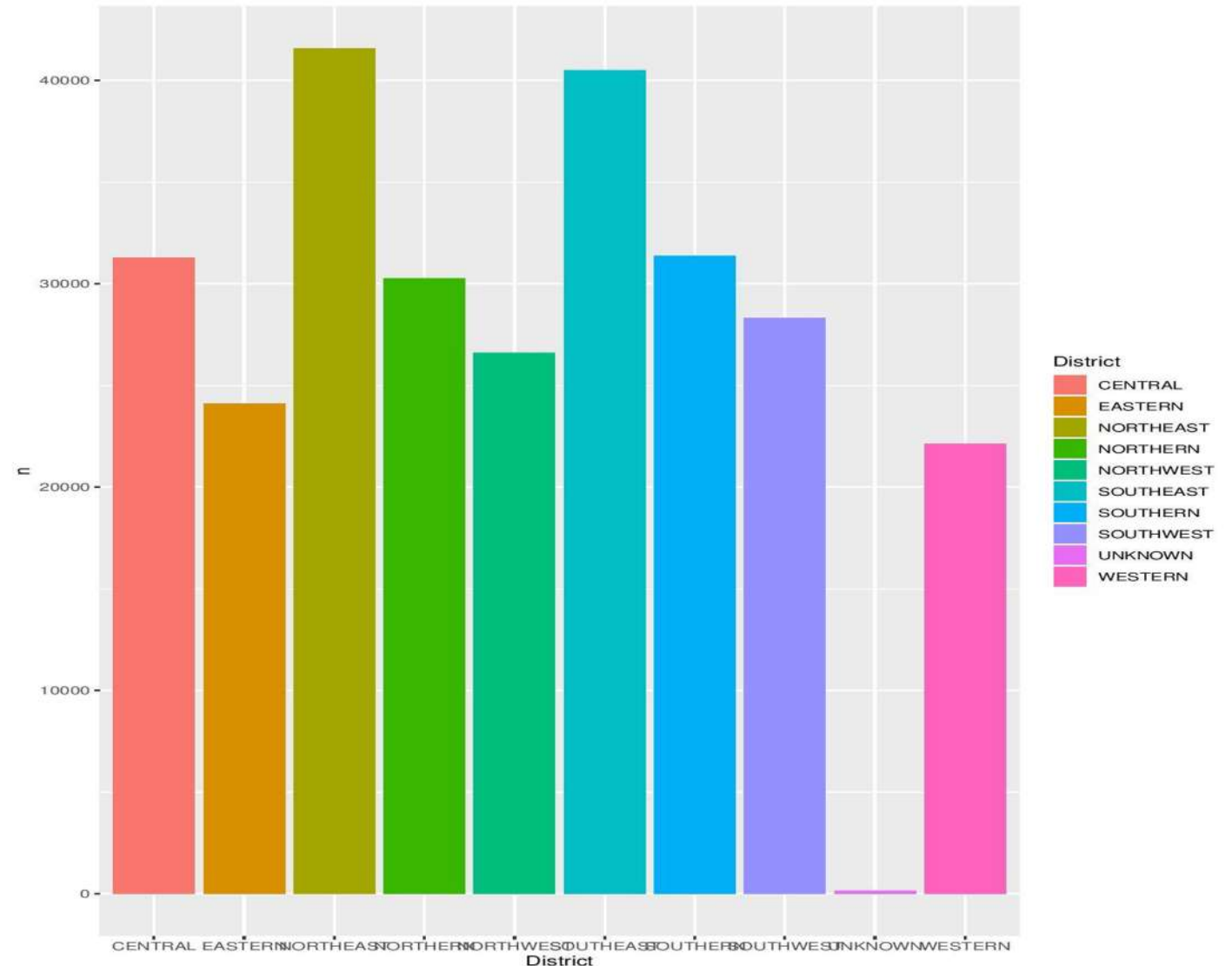
## Evolution of crimes based on the District

### Hypothesis:

- Correlation between the area and the number of crimes

### Results:

- Number of crimes vary amongst the districts
- District with the most crimes: Northeast
- Less crimes within the Western and Eastern



# ● Data Analysis

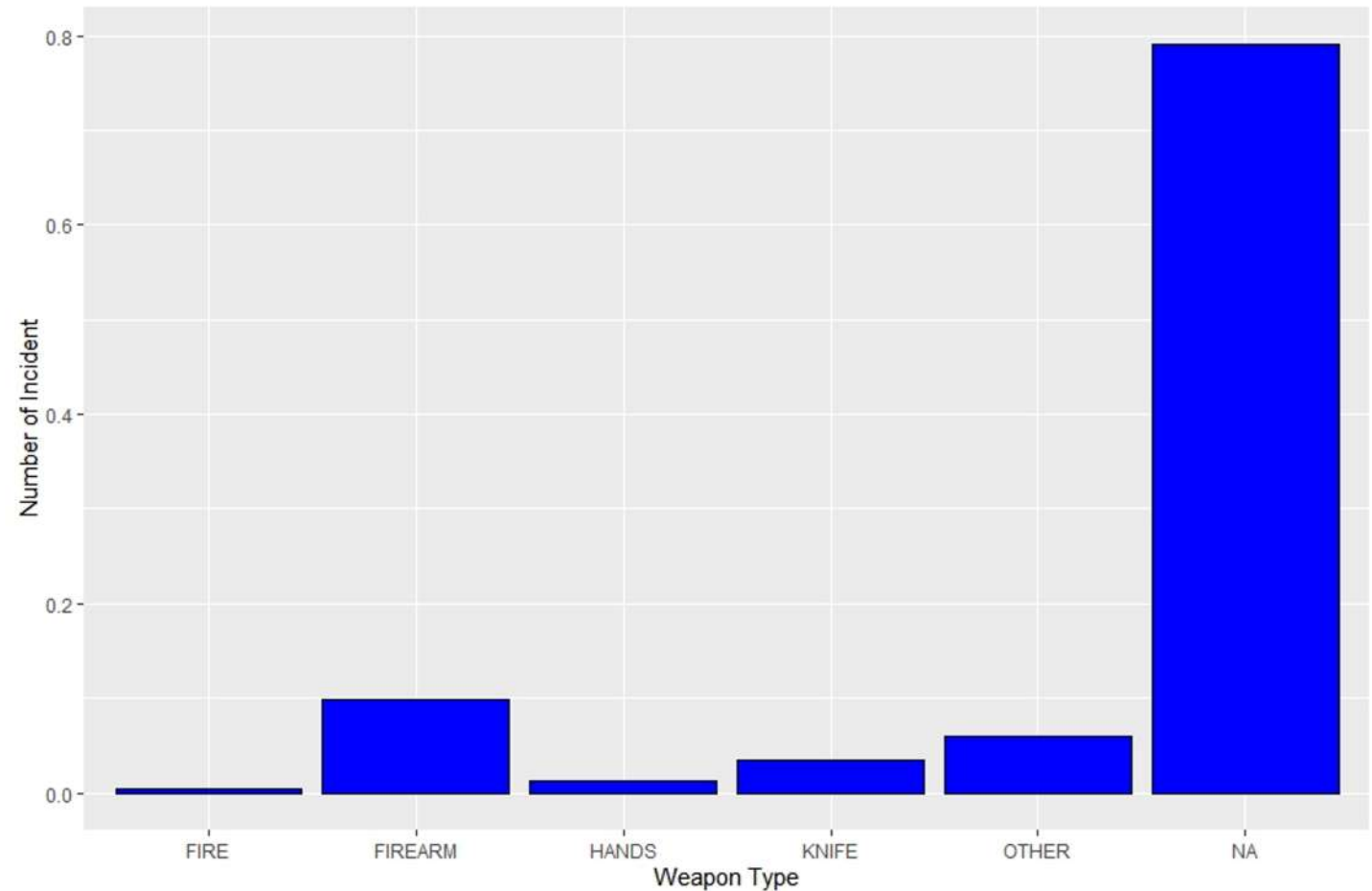
## Classification of Crimes Based on Weapon Used

### Hypothesis:

- Classifying the number of crimes based on weapon used

### Results:

- Most crimes committed without any weapon
- Subsequently, 'Firearm' dominates the crime count



# ● Data Analysis

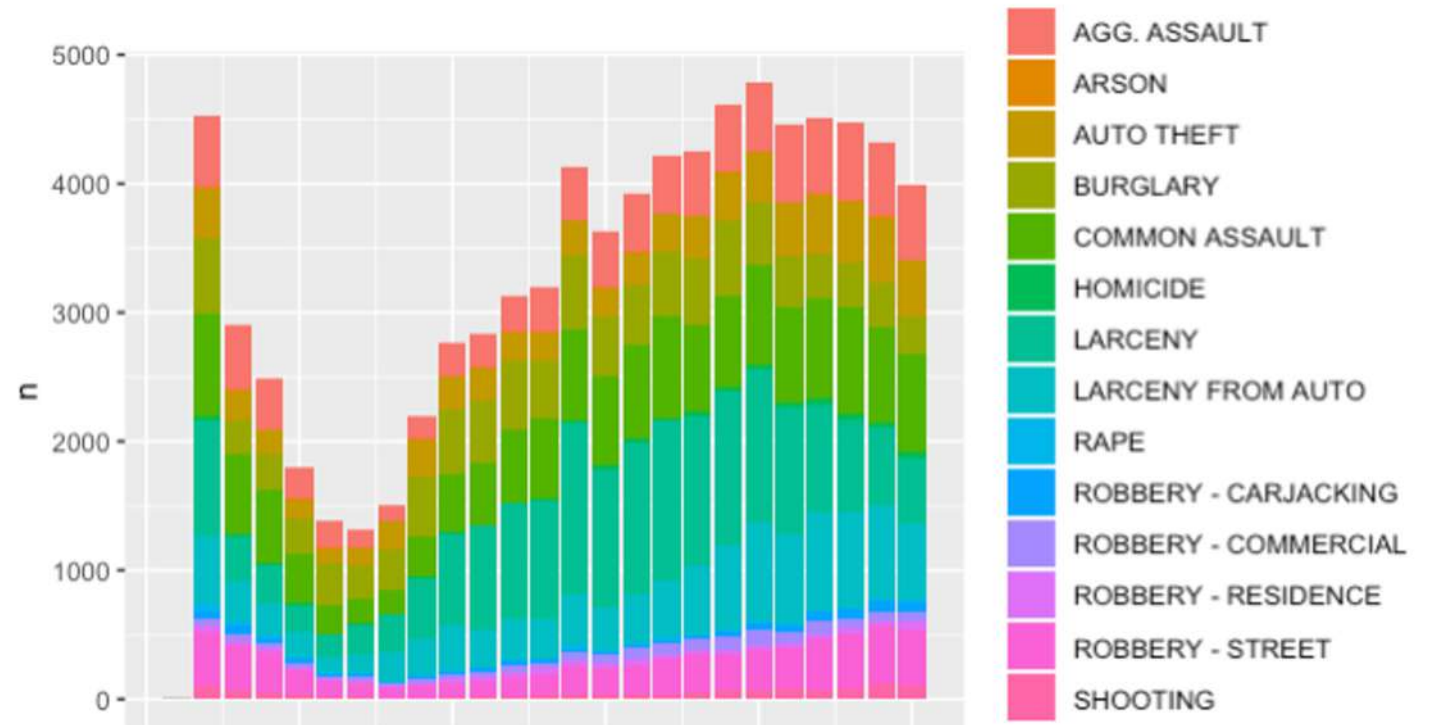
Evolution of crimes based on the time

## Hypothesis:

- Correlation between the respective hour of the day and the types of crimes

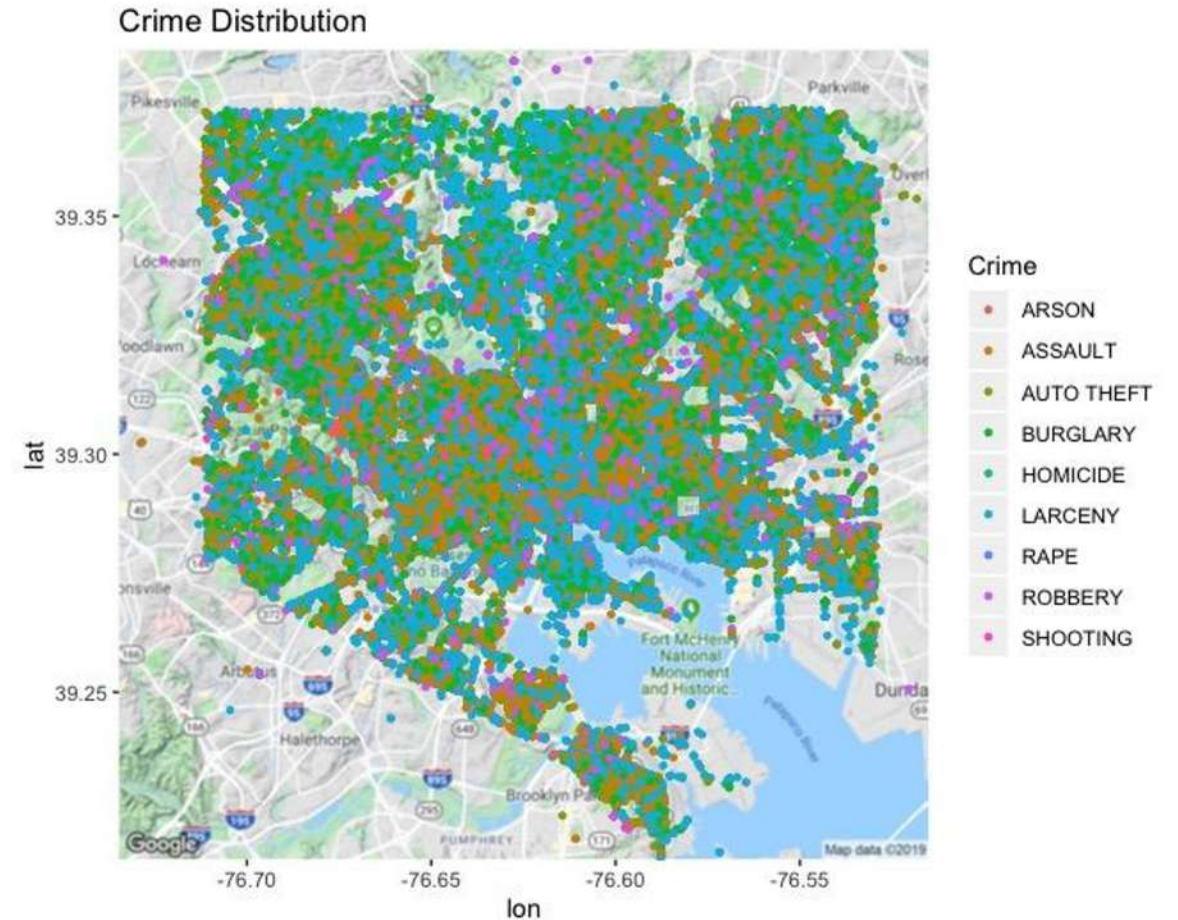
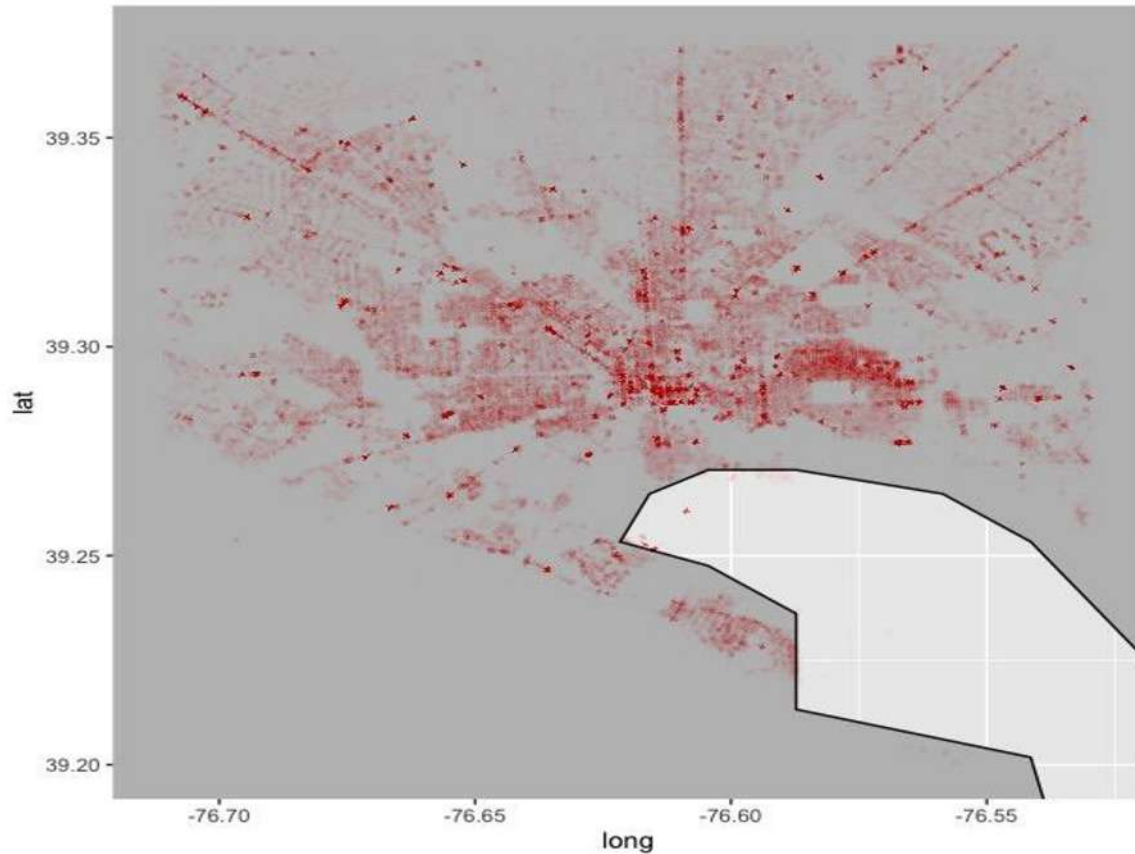
## Results:

- Data indicates that crime increases throughout the day.
- Crime spikes after 09:00.
- Crime is at it's highest at 20:00.



# ● Data Analysis

## Crime distribution



Graphical representation of crime throughout Baltimore City



# ● Data Analysis

```
72 #DATA ANALYSIS
73
74 #Data Cleaning
75 #Subsetting rows containing missing values
76 crime_dat_miss <- subset(crime_dat, select = c(CrimeDate, CrimeTime, Location, Description, District, Longitude, Latitude, year, month))
77 crime_dat_miss <- crime_dat_miss[complete.cases(crime_dat_miss), ]
78
79 crime_dat_miss$hour <- substring(crime_dat_miss$CrimeTime, 1, 2)
80 #crime_dat_miss$TimeOfDay <- with(crime_dat_miss, ifelse(hour >= 5 & hour < 12, "Morning",
81 #                                     ifelse(hour >= 12 & hour <= 17, "Afternoon", ifelse(hour >= 17 & hour <= 22, "Evening", "Night"))))
82 crime_dat_miss$TimeOfDay <- with(crime_dat_miss, ifelse(hour >= 5 & hour < 6, "Day", "Night"))
83
84 crime_dat_miss$hour <- as.factor(crime_dat_miss$hour)
85
86 crime_dat_miss$CrimeDate <- as.Date(crime_dat_miss$CrimeDate, format = "%m/%d/%Y")
87 crime_dat_miss$day <- weekdays(as.Date(crime_dat_miss$CrimeDate))
88
89
90 crime_dat_miss$DayOfWeek <- with(crime_dat_miss, ifelse(day == "Monday", "1",
91                                                         ifelse(day == "Tuesday", "2",
92                                                         ifelse(day == "Wednesday", "3",
93                                                         ifelse(day == "Thursday", "4",
94                                                         ifelse(day == "Friday", "5",
95                                                         ifelse(day == "Saturday", "6", "7"))))))))
96
97 #crime_dat_miss$HourOfDay <-
98 #crime_dat_miss$Location <- as.character(crime_dat_miss$Location)
99
```

# ● Data Analysis

```
23 install.packages("cluster")
24 library(cluster)
25 crime_dat<-read.csv("BCD.csv",header = TRUE, sep = ",")
26
27 #Extracting and making new columns for year and month
28 crime_dat$year<-substring(crime_dat$CrimeDate,7,11)
29 crime_dat$month<-substring(crime_dat$CrimeDate,1,2)
30 temporary_dataset<-crime_dat[,c("year","CrimeCode","Description")]
31
32
33 #ddply(crime_dat, "crime_dat", summarize)
34
35 #Count of crimes by type of crime
36 crime_count<-temporary_dataset %>% group_by(Description) %>% tally() %>% arrange(desc(n))
37
38 #unique(crime_dat$year)
39 #Count of crimes by type of crime, yearwise
40 crime_year <- temporary_dataset %>% group_by(year,Description) %>% filter(year > "2013") %>% tally()
41 ggplot(crime_year,aes(year,n,colour = Description, group=1))+geom_point() +xlab("Year")+ylab("Count")
42
43 #Count of crimes by year
44 crime_count_year <- temporary_dataset %>% group_by(year) %>% filter(year > "2013") %>% tally() %>% arrange(desc(n))
45 ggplot(crime_count_year,aes(year,n,fill = year)) + geom_bar(stat="identity")
46
47
48 unique(crime_dat$District)
49 #Bar Plot of description and count per year
50 ggplot(crime_count,aes(Description,n)) + geom_bar(stat="identity")
51
52 #Analyzing by district
53 crime_district<-crime_dat[,c("District","year","CrimeCode","Description")] %>% group_by(District)
54
55 #Crime count by district
56 crime_count_district<-crime_district %>% group_by(District) %>% tally() %>% arrange(desc(n))
57 ggplot(crime_count_district,aes(District,n,fill=District)) + geom_bar(stat="identity")
58
59 #Crime count by month
60 temporary_dataset<-crime_dat[,c("month","CrimeCode","Description")]
61 crime_month <- temporary_dataset %>% group_by(month) %>% tally()
62 ggplot(crime_month,aes(month,n,fill=month)) + geom_bar(stat="identity")
63
```

# ● Data Analysis

```
23 install.packages("cluster")
24 library(cluster)
25 crime_dat<-read.csv("BCD.csv",header = TRUE, sep = ",")
26
27 #Extracting and making new columns for year and month
28 crime_dat$year<-substring(crime_dat$CrimeDate,7,11)
29 crime_dat$month<-substring(crime_dat$CrimeDate,1,2)
30 temporary_dataset<-crime_dat[,c("year","CrimeCode","Description")]
31
32
33 #ddply(crime_dat, "crime_dat", summarize)
34
35 #Count of crimes by type of crime
36 crime_count<-temporary_dataset %>% group_by(Description) %>% tally() %>% arrange(desc(n))
37
38 #unique(crime_dat$year)
39 #Count of crimes by type of crime, yearwise
40 crime_year <- temporary_dataset %>% group_by(year,Description) %>% filter(year > "2013") %>% tally()
41 ggplot(crime_year,aes(year,n,colour = Description, group=1))+geom_point() +xlab("Year")+ylab("Count")
42
43 #Count of crimes by year
44 crime_count_year <- temporary_dataset %>% group_by(year) %>% filter(year > "2013") %>% tally() %>% arrange(desc(n))
45 ggplot(crime_count_year,aes(year,n,fill = year)) + geom_bar(stat="identity")
46
47
48 unique(crime_dat$District)
49 #Bar Plot of description and count per year
50 ggplot(crime_count,aes(Description,n)) + geom_bar(stat="identity")
51
52 #Analyzing by district
53 crime_district<-crime_dat[,c("District","year","CrimeCode","Description")] %>% group_by(District)
54
55 #Crime count by district
56 crime_count_district<-crime_district %>% group_by(District) %>% tally() %>% arrange(desc(n))
57 ggplot(crime_count_district,aes(District,n,fill=District)) + geom_bar(stat="identity")
58
59 #Crime count by month
60 temporary_dataset<-crime_dat[,c("month","CrimeCode","Description")]
61 crime_month <- temporary_dataset %>% group_by(month) %>% tally()
62 ggplot(crime_month,aes(month,n,fill=month)) + geom_bar(stat="identity")
63
```



# ● Data Analysis

```
116 #K Means Clustering
117 install.packages("fpc")
118 library(fpc)
119
120 set.seed(20)
121 CrimeCluster <- kmeans(crime_dat_miss[, 4:5], 9, nstart = 20)
122 crime_dat_miss$loc <- as.factor(CrimeCluster$cluster)
123 str(CrimeCluster)
124
125 library(ggmap)
126
127 #ggmap
128 BaltimoreMap <- get_map("Maryland", zoom = 10)
129 ggmap(BaltimoreMap) + geom_point(aes(x = Longitude, y = Latitude, colour = as.factor(crime_dat_miss$loc)), data = crime_dat_miss) +
130   ggtitle("NCrimes using KMean")
131
132 ggplot(data = subset(states, region == c("maryland")), mapping = aes(x = long, y = lat, group = group)) + coord_fixed(1.3) +
133   geom_polygon(color = "black", fill = "gray") +
134   geom_point(data=crime_dat_miss, aes(x = crime_dat_miss$Longitude, y = crime_dat_miss$Latitude, group=crime_dat_miss$loc), pch=21, size=1, alpha=I(0.005)) +
135   coord_fixed(xlim = c(-76.71162, -76.5285), ylim = c(39.20041, 39.37293) )
136
137 register_google(key = "AIzaSyDjLgqH4IJJdmQyBIVHKlB7eaacFYS9fV4", write = TRUE)
138
139
140 clusplot(CrimeCluster, CrimeCluster$cluster, color=TRUE, shade=TRUE, labels=2, lines=0)
141 plotcluster(CrimeCluster, CrimeCluster$cluster)
142
```

# ● Data Analysis

```
143 #Decision Tree
144 set.seed(1000)
145 #Creating training and testing datasets
146
147 crime_dat_miss <- crime_dat_miss %>% filter(year > "2017")
148 crime_dat_miss$hour <- as.numeric(crime_dat_miss$hour)
149 crime_dat_miss$TimeOfDay <- with(crime_dat_miss, ifelse(hour >= 5 & hour<12, "Morning",
150 ifelse(hour>=12 & hour<=17, "Afternoon", ifelse(hour>=17 & hour<=20, "Evening", "Night"))))
151 crime_dat_miss$DayOfWeek <- as.factor(crime_dat_miss$DayOfWeek)
152 ind <- sample.split(Y=crime_dat_miss$CrimeCode, SplitRatio = 0.75)
153 crime_dat_miss$CrimeCode <- as.numeric(crime_dat_miss$CrimeCode)
154 train_set <- crime_dat_miss[ind,]
155 test_set <- crime_dat_miss[!ind,]
156
157 CrimeTree <- rpart(District ~ TimeOfDay + DayOfWeek + hour + Description, data = train_set , method = "class")
158 par(xpd = NA) # Avoid clipping the text in some device
159 plot(CrimeTree)
160 text(CrimeTree)
161 predicted.classes <- CrimeTree %>% predict(test_set, type ="class")
162 mean(predicted.classes == test_set$CrimeCode)
163
164 #KNN Classification
165 normalize <- function(x) {
166   return ((x - min(x)) / (max(x) - min(x))) }
167 prc_n <- as.data.frame(lapply(crime_dat_miss[4:5], normalize))
168 train_set <- crime_dat_miss[ind,]
169 test_set <- crime_dat_miss[!ind,]
170 prc_test_pred <- knn(train = train_set, test = test_set,cl = train_set$Crime, k=10)
171
172
```



## ● Data Analysis

```
~/Documents/Academics/Term 1/Big Data/R Project/
70470 11153 0.86335959
model <- randomForest(CrimeViolent ~ DayOfWeek + Season + Lat
iss, ntree=200)
model

Call:
randomForest(formula = CrimeViolent ~ DayOfWeek + Season + Lat
ne_dat_miss, ntree = 200)
Type of random forest: classification
Number of trees: 200
No. of variables tried at each split: 2

OOB estimate of error rate: 25.12%
Confusion matrix:
      0      1 class.error
0 178329 14739 0.07634098
1  54271 27352 0.66489837
>
```

# ● Conclusions

## Final Remarks



Clear indication from the data used that crime and location correlate, as well as certain measures of time.



Use the latitude and longitude in order to map the crimes □ Correlation between road axes and crimes



# THANK YOU

