



Analysis of Crimes in the City of New York

Capstone Project

Michał Czapski
czapski.michal@gmail.com

Contents

- **Introduction**
- **Data Exploration**
- **Data Modeling**
 - Supervised Learning
 - Unsupervised Learning
- **Assumptions and Limitations**
- **Recommendations**

Introduction

New York City is the most populous city in the USA
divided into 5 different boroughs



Every borough differs in demographics, wealth and lifestyle



Manhattan



Bronx



Queens



Staten Island



Brooklyn

The goal of this study is to investigate crime rates in NYC

The main focus is to investigate:

- Changes in crime rates between 2006 and 2016
- Differences between boroughs in crime rates
- Crime rates and demographics / housing market relationships



There are three potential clients that can be interested in this study



NYPD



Housing Market



NYC Authorities

Three main datasets were used in this analysis

NYPD Complaint Data

contains information about type of offense, time of occurrence, specific location and borough.

NYU Furman Center

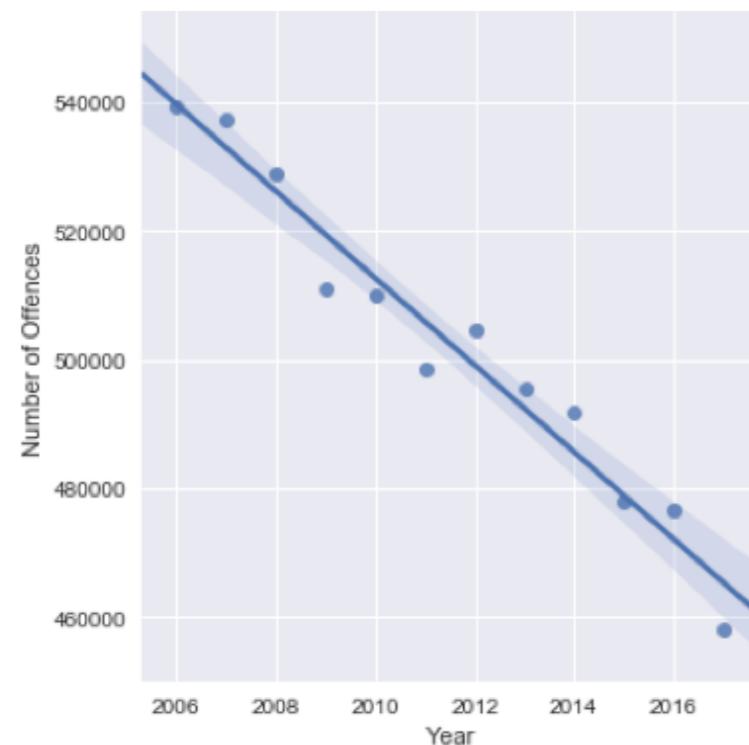
provides a lot of information on neighborhood demographics between 2010 and 2015

NYC Dept. of Finance

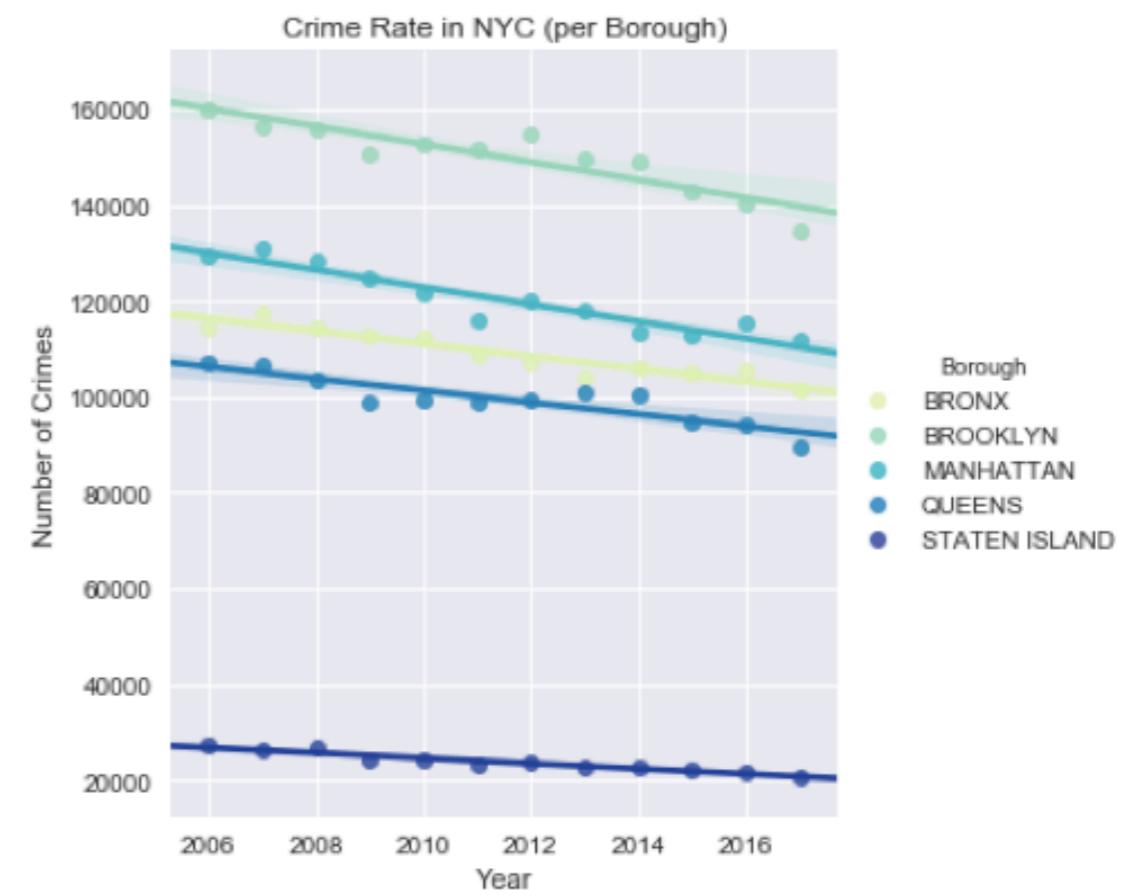
provides information on property prices (means, medians) and number of sales per year per borough per type of home.

Data Exploration

On average the number of reported crimes decline every year on the city and borough level

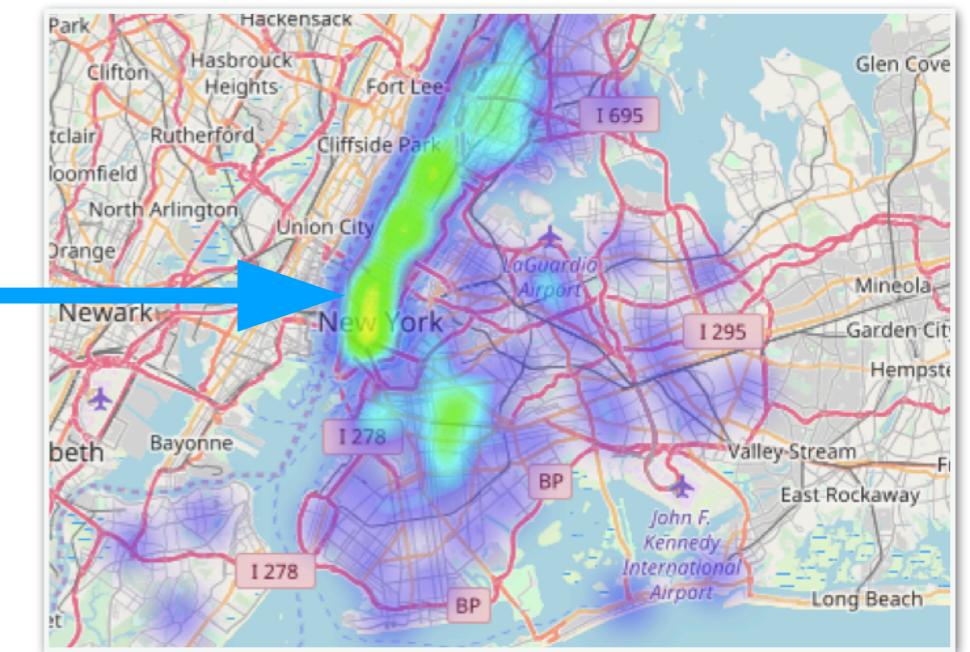
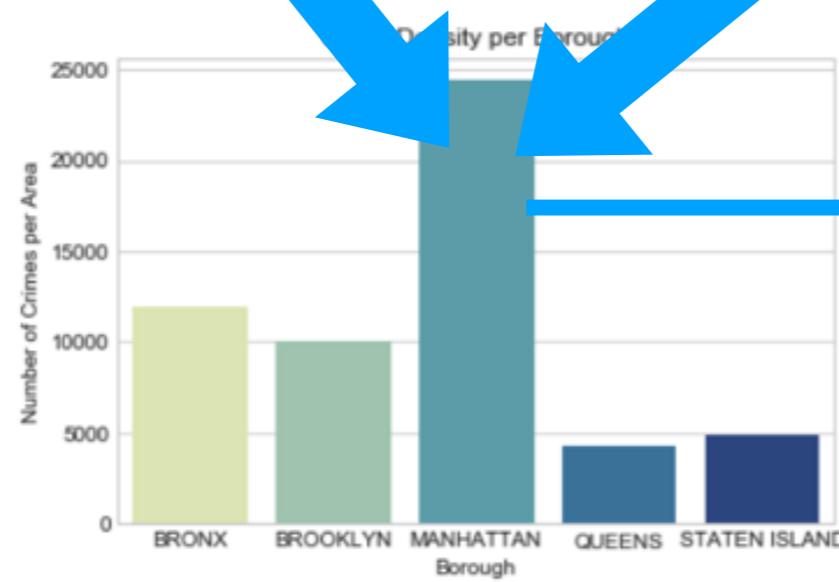
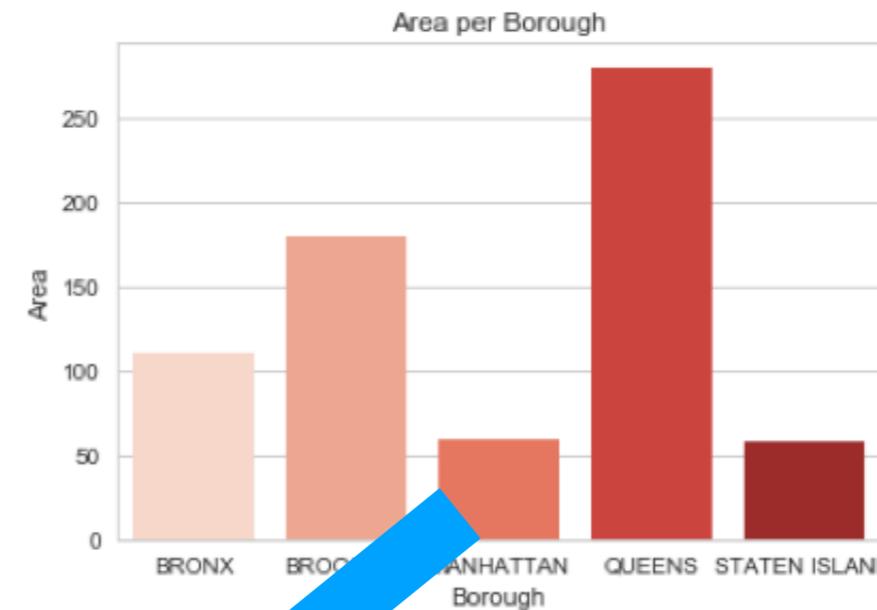
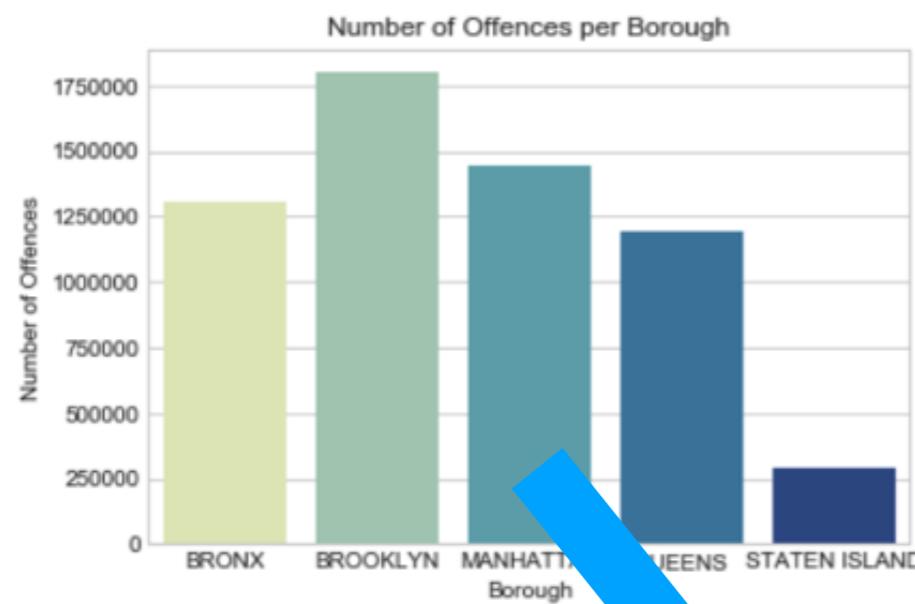


NYC
-7000 crimes / year



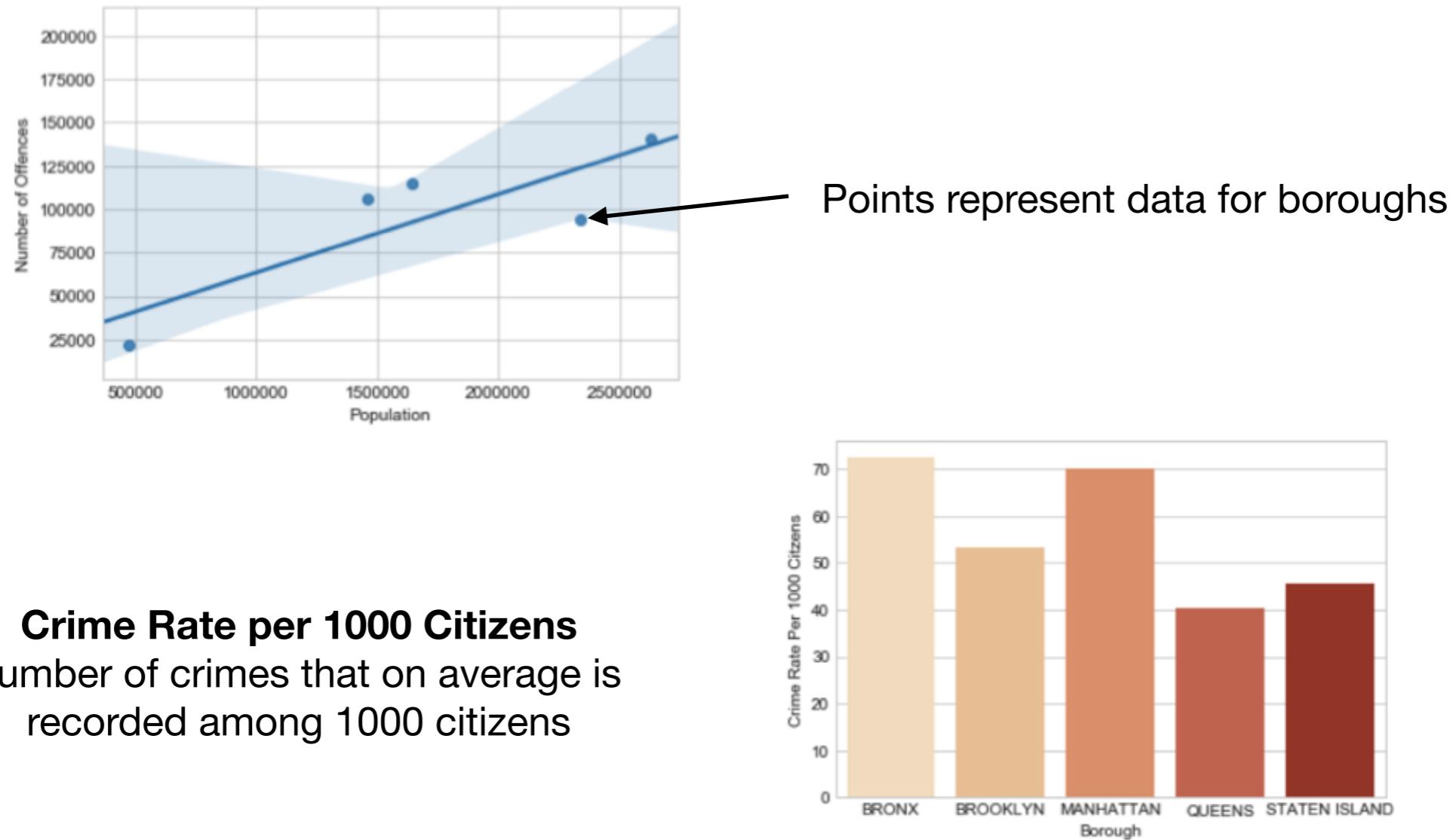
max (Brooklyn)
-2000 crimes per year
min (Staten Island)
-550 crimes per year

Boroughs differ significantly in areas, (reported) crime density can better compare borough-level crime rates

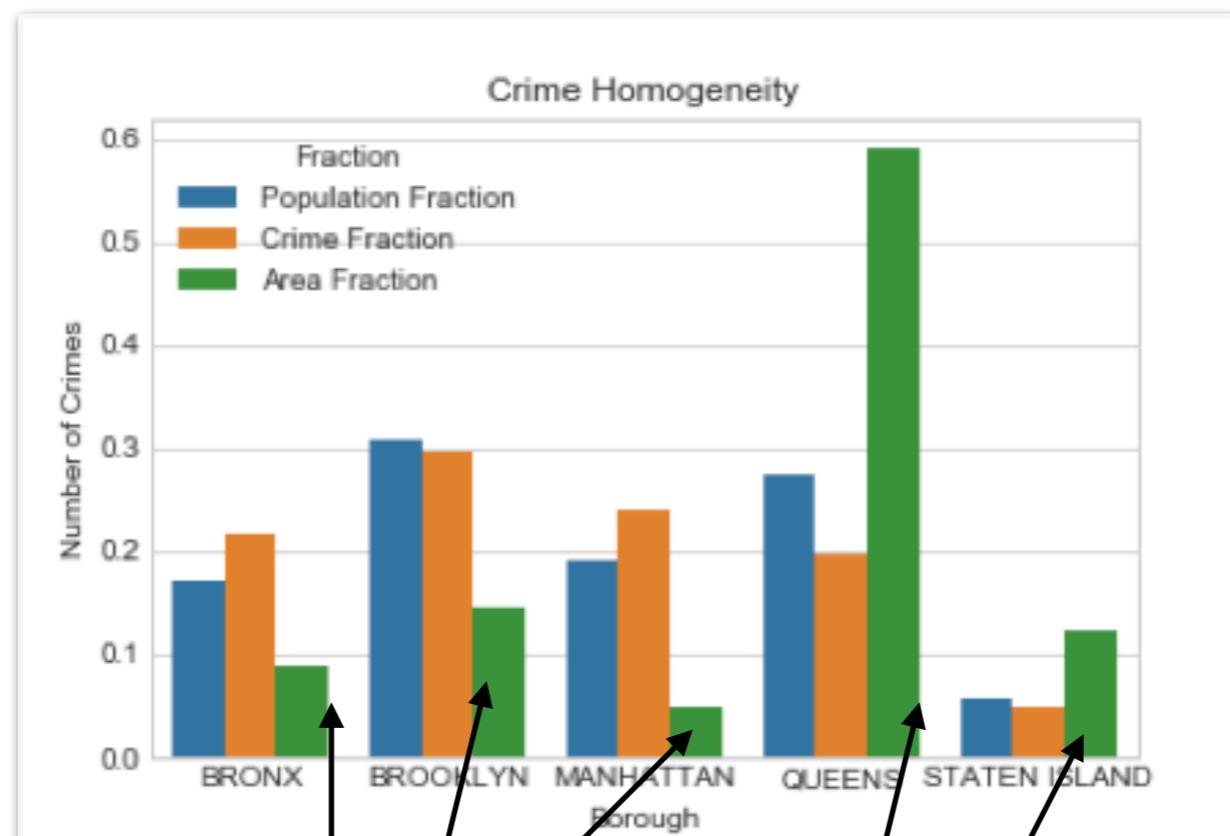


Manhattan outdistanced other boroughs

Population also can impact crime rates, the bigger population generally the likelihood of crime should be generally greater



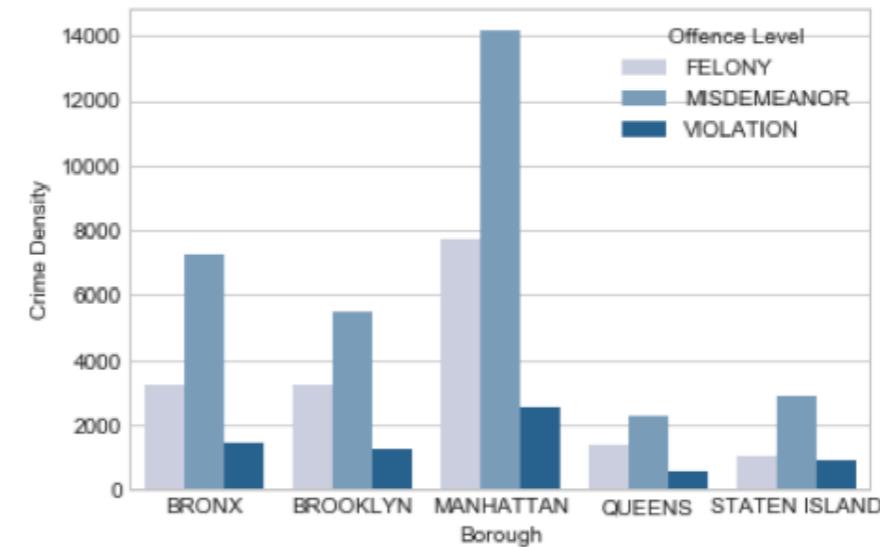
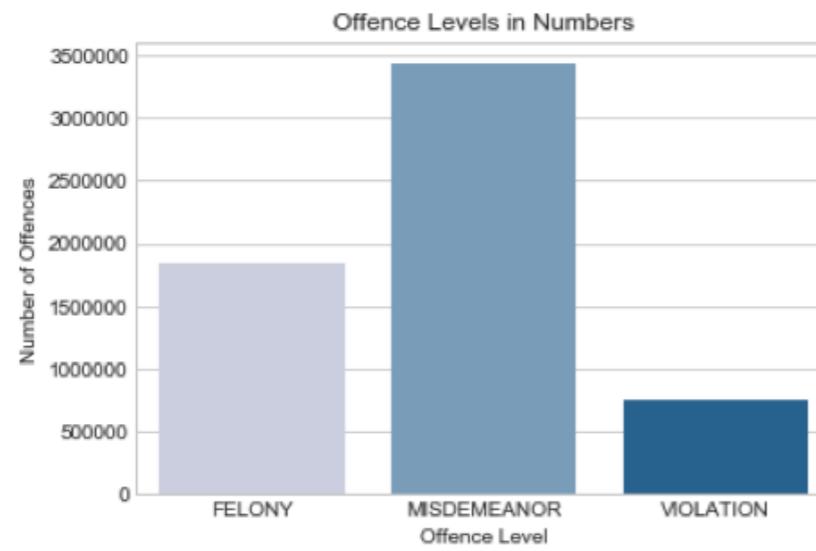
There are disproportions in boroughs' homogeneity in terms of population, area and crime rates share



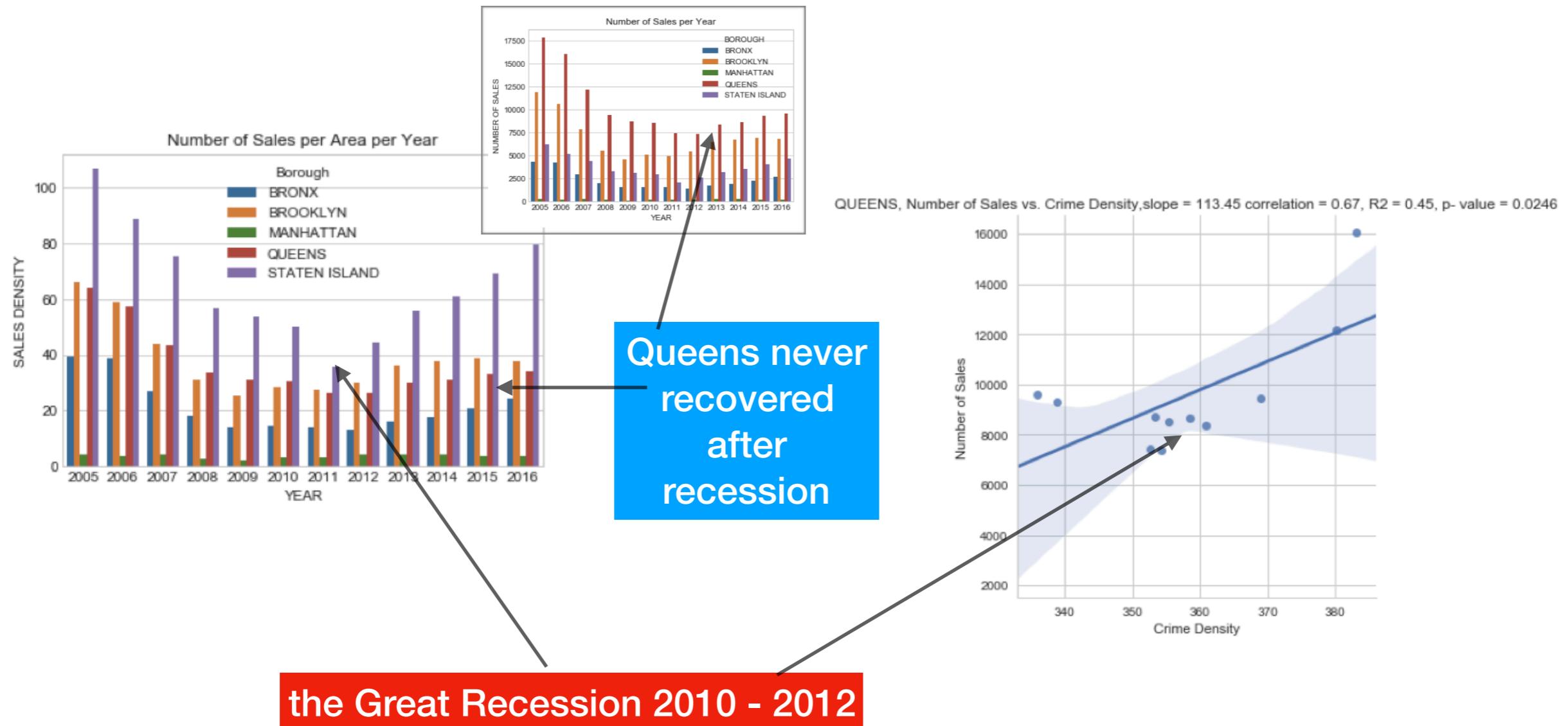
Area is the smallest fraction bigger fraction

Area is the largest fraction bigger fraction

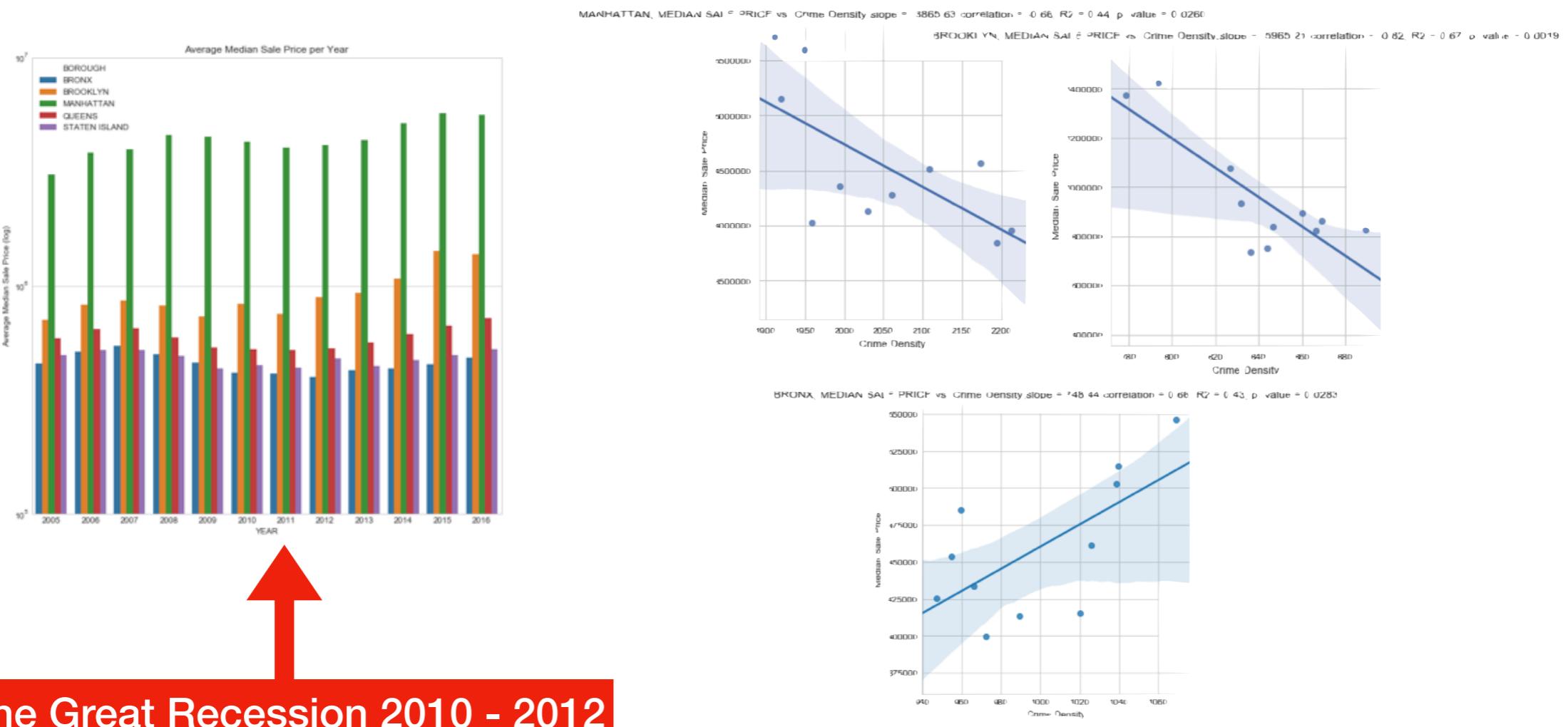
Generally distribution of different crime level is similar in all three boroughs i.e. misdemeanors are twice more frequent than felonies and five times more frequent than violations



Sales density is a better indicator to compare number of sales, the Great Recession had significant impact on sales



Average median sale prices also were affected by the Recession after which Manhattan and Brooklyn prices went up and in the Bronx they never recovered

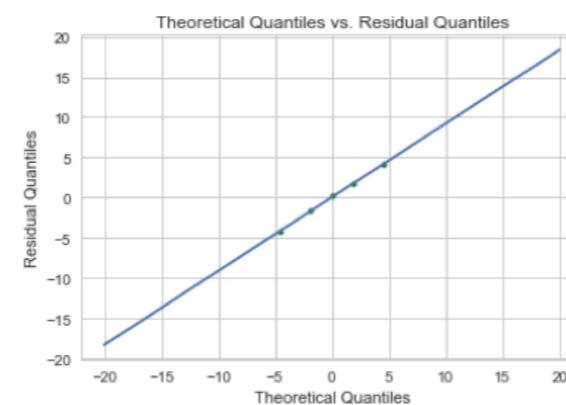
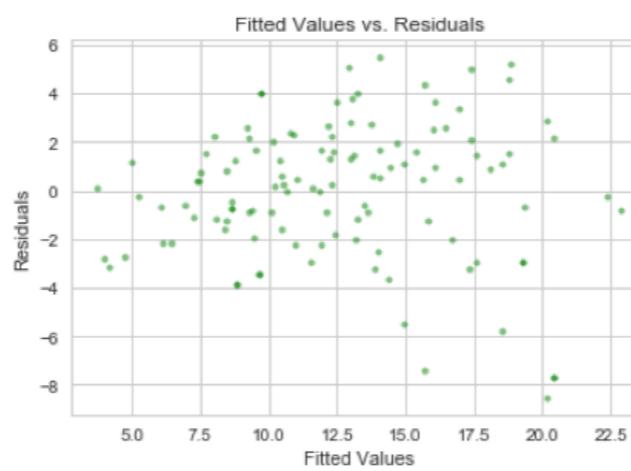
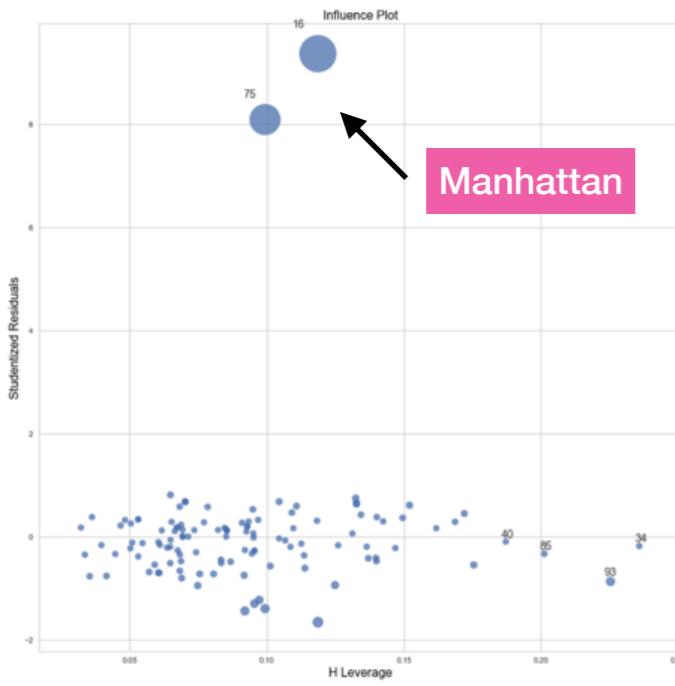


Data Modeling

Supervised Learning

Regression model built is not the most efficient in predicting values, but it can help to understand which demographic factors influence sever crime rates (per 1000 citizens)

Response Variable:
Severe Crime Rates
(per 1000 citizens)



Variable	Coefficient
Intercept	-598.95
Year	0.33
Born in New York State	-41.78
Income diversity ratio	-38.73
Percent Hispanic	-4.29
Foreign-born population	-19.34
Percent Asian	-12.22
Population aged 65+	-84.48
Single-person households	22.76
Population density (1,000 persons per square mile)	-0.08

Cross validation average score	0.5077
R-squared (R2)	0.6
Residual standard error (RSE)	8.48
Average response	12.6
RSE / Average response	0.67

Unsupervised Learning

The goal of unsupervised learning is to check how strong are the similarities between the 77 NYC precincts and if there are any homogenous groupings among them

Features: Offence Codes

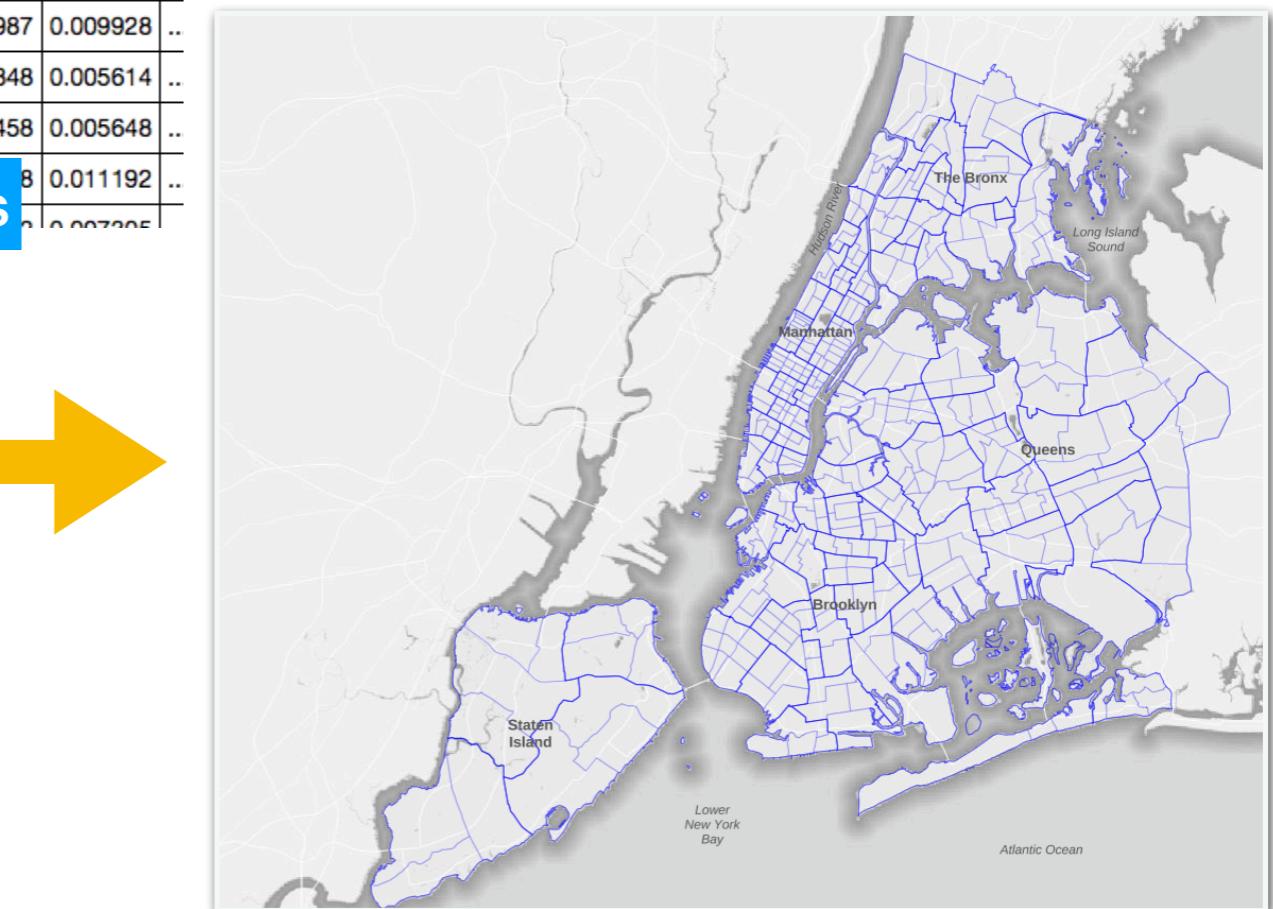


Offence Code	101	102	103	104	105	106	107	109	110	111	..
0	0.000558	0.000105	0.000000	0.003695	0.037196	0.036429	0.070731	0.467960	0.018755	0.005822	..
1	0.000937	0.000117	0.000000	0.004861	0.079656	0.107477	0.082116	0.396112	0.015521	0.005798	..
2	0.000554	0.000000	0.000000	0.004985	0.077550	0.058674	0.092889	0.570801	0.021987	0.009928	..
3	0.001295	0.000000	0.000072	0.008780	0.127600	0.111407	0.073983	0.308672	0.036848	0.005614	..
4	0.001031	0.000000	0.000000	0.007256	0.085836	0.077384	0.095442	0.411328	0.027458	0.005648	..
5	0.000933	0.000124	0.000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.011192	..
6	0.000558	0.000000	0.000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.007205	..

Offence frequency vectors

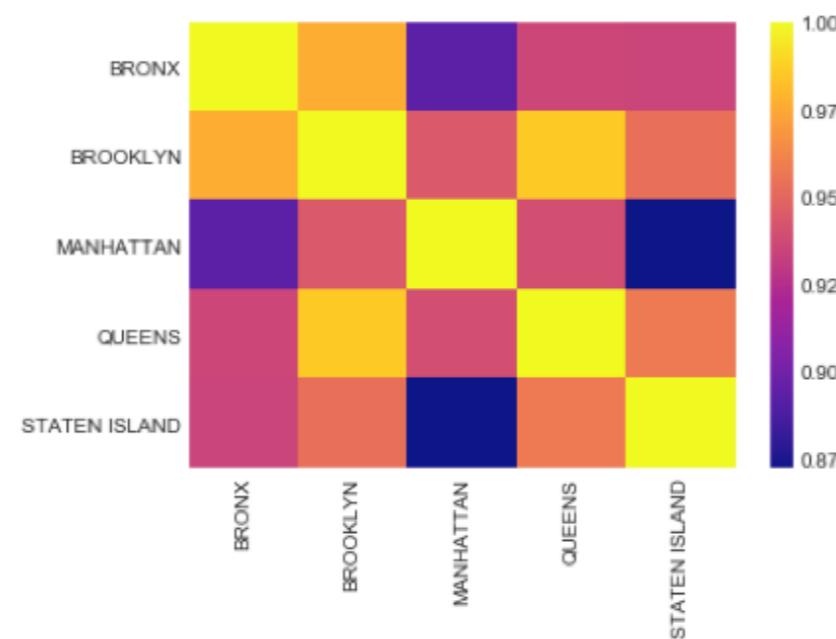


Observations: Precincts



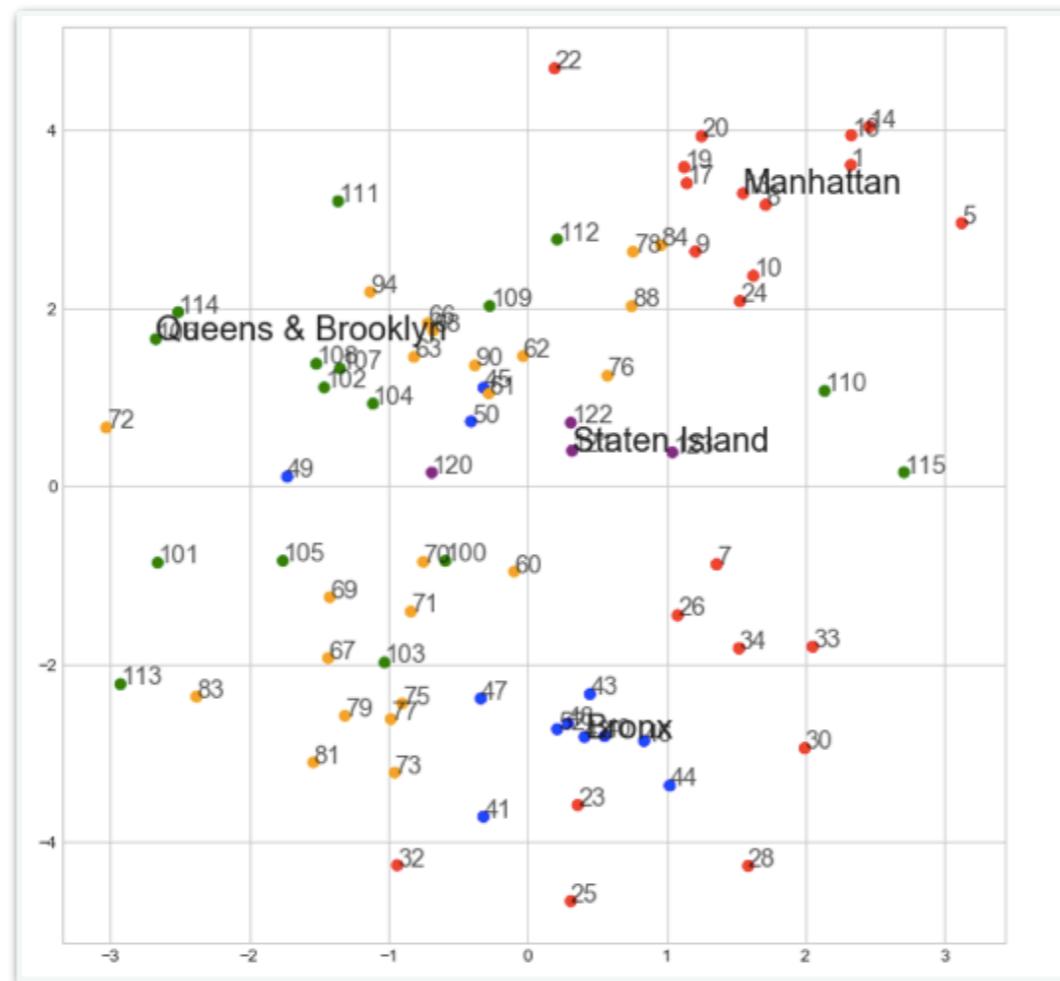
Boroughs generally are very similar to each other crime-pattern wise with Manhattan being the most different from all other boroughs and Brooklyn and Queens showing the closest resemblance

	BRONX	BROOKLYN	MANHATTAN	QUEENS	STATEN ISLAND
BRONX	1.000	0.977	0.894	0.936	0.935
BROOKLYN	0.977	1.000	0.944	0.986	0.954
MANHATTAN	0.894	0.944	1.000	0.940	0.872
QUEENS	0.936	0.986	0.940	1.000	0.958
STATEN ISLAND	0.935	0.954	0.872	0.958	1.000

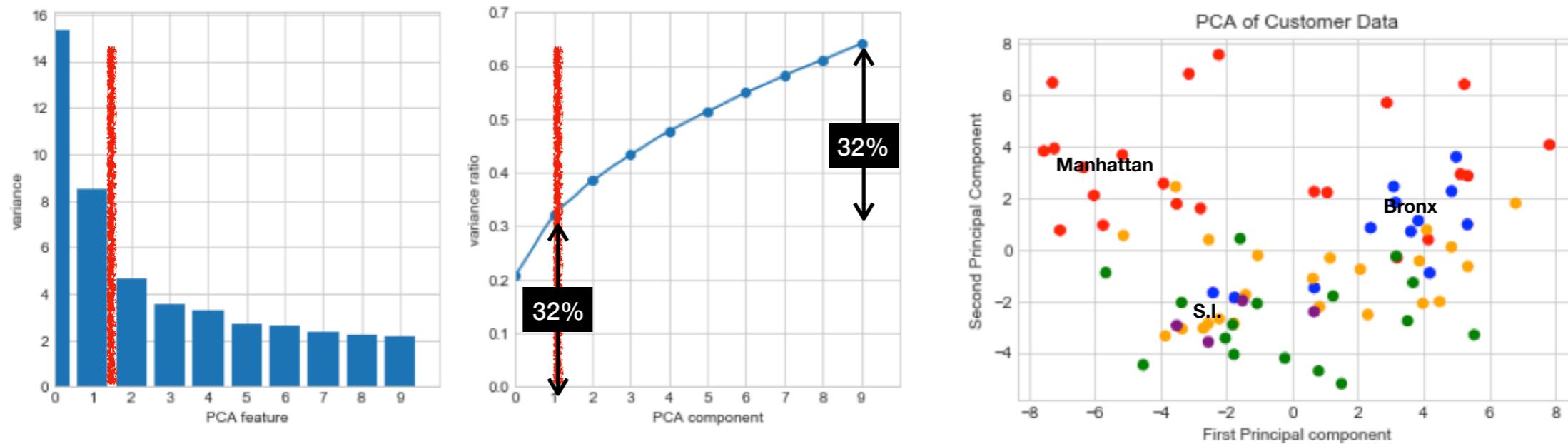


T-SNE decomposition used to transform offence code frequencies into two dimensional space, showed some visual clustering of the precincts

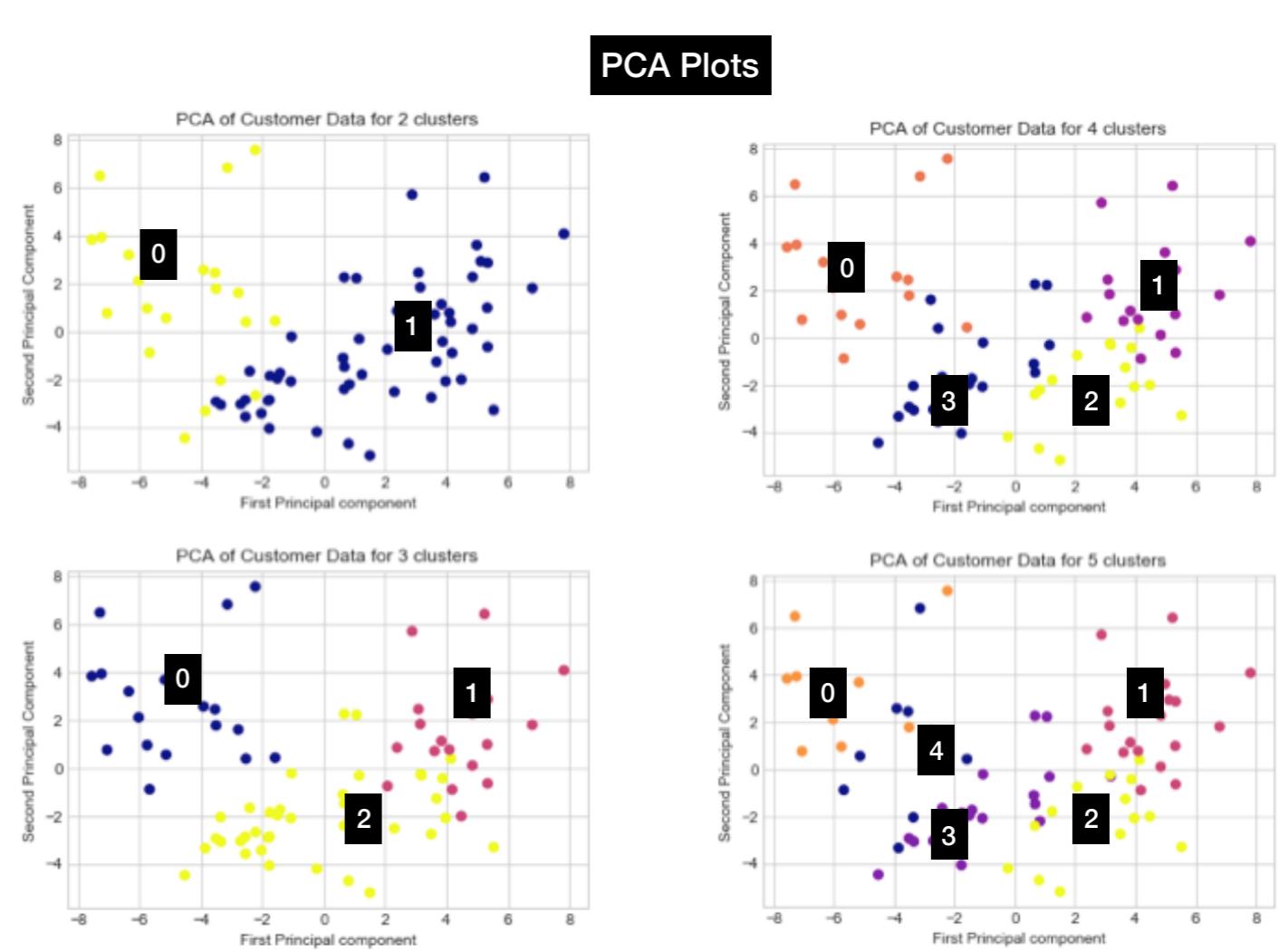
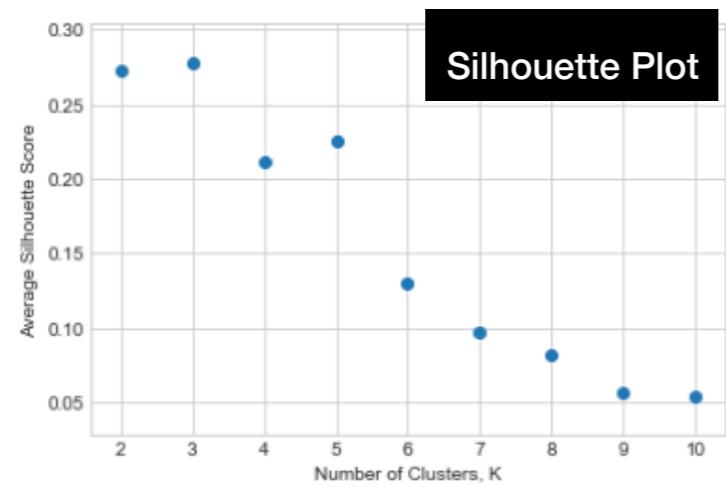
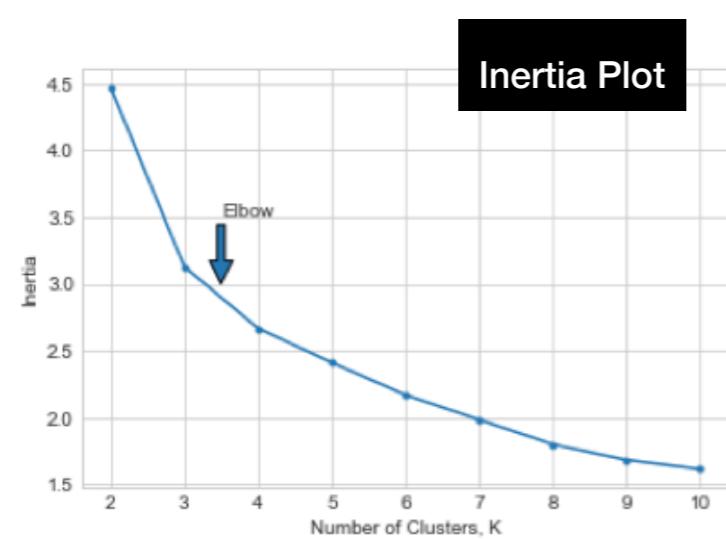
t-SNE preserves
nearness of the
samples



PCA analysis of variance showed that the dataset can be decomposed into two intrinsic dimensions, that revealed similar clustering to t-SNE method



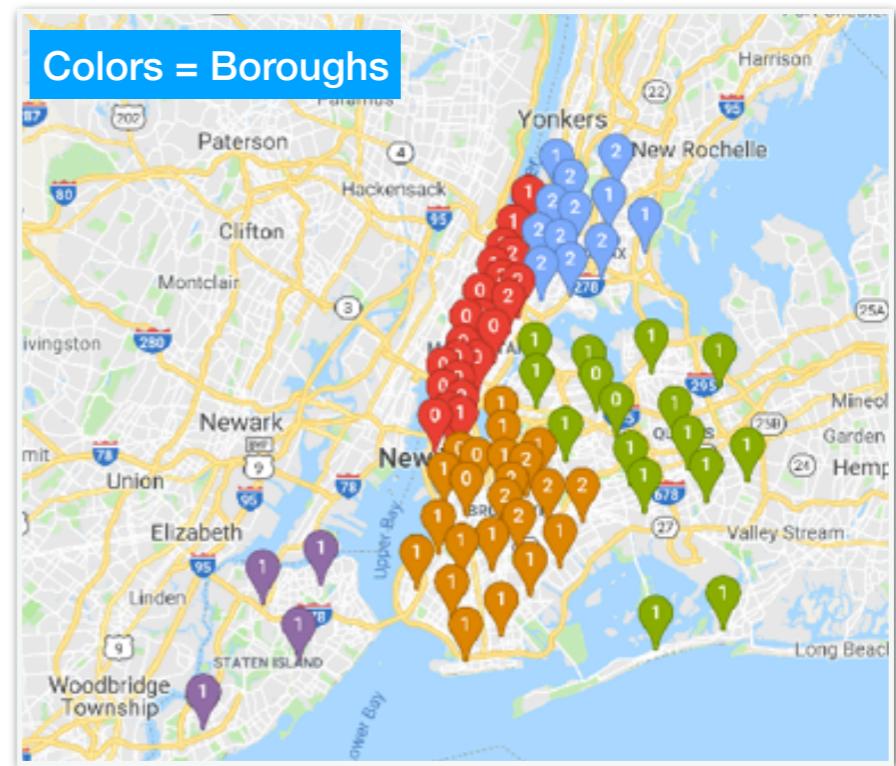
Three methods that were used to specify the number of clusters for K-Means clustering method, give similar results of best K around 2 - 3 clusters



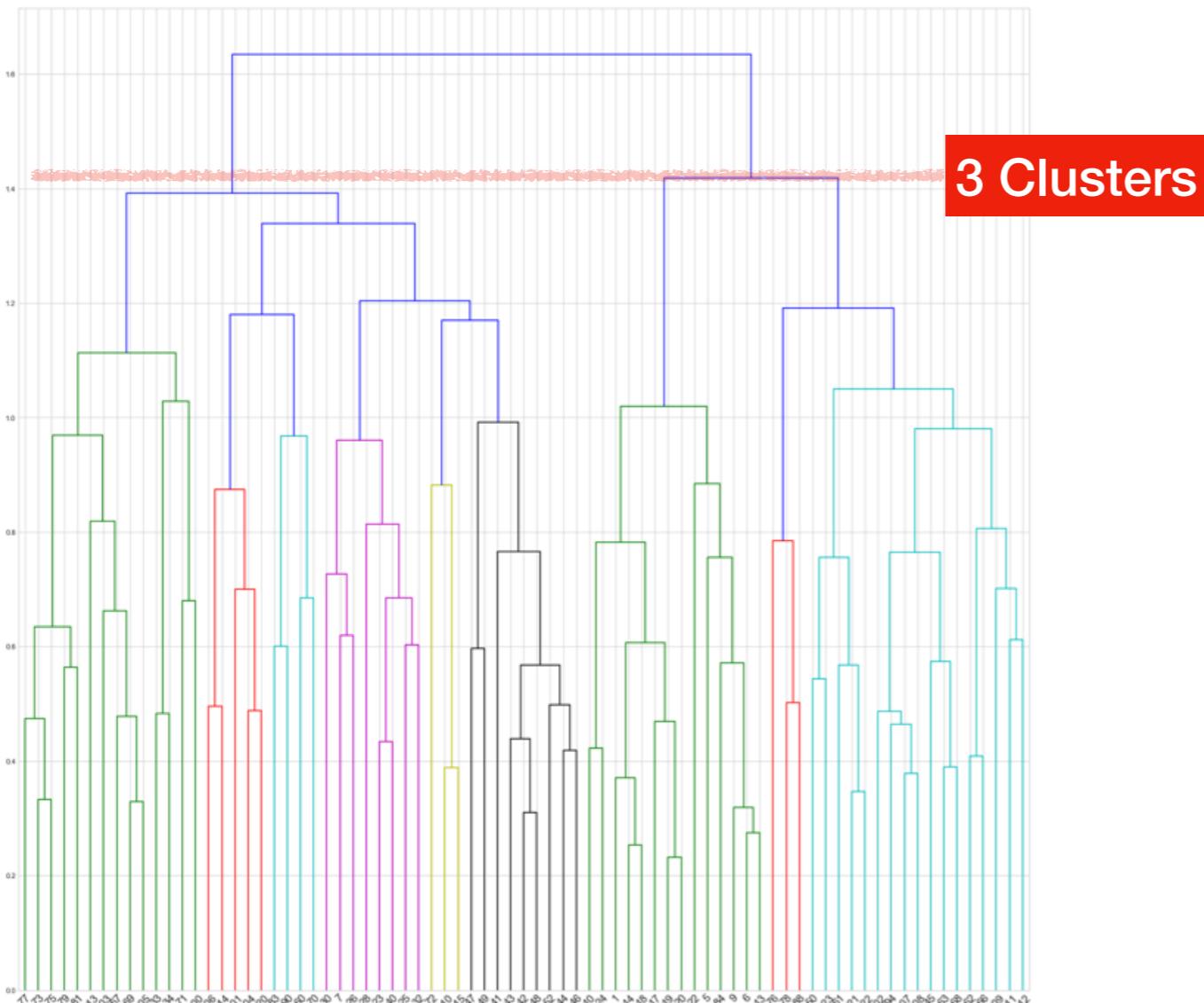
K-Means method for K = 3 produced clusters that often group neighboring precincts. All clusters show similar most frequent offence codes

LABEL	BRONX	BROOKLYN	MANHATTAN	QUEENS	STATEN ISLAND
0	0	0	8	1	0
1	5	6	5	8	0
2	0	7	0	2	1

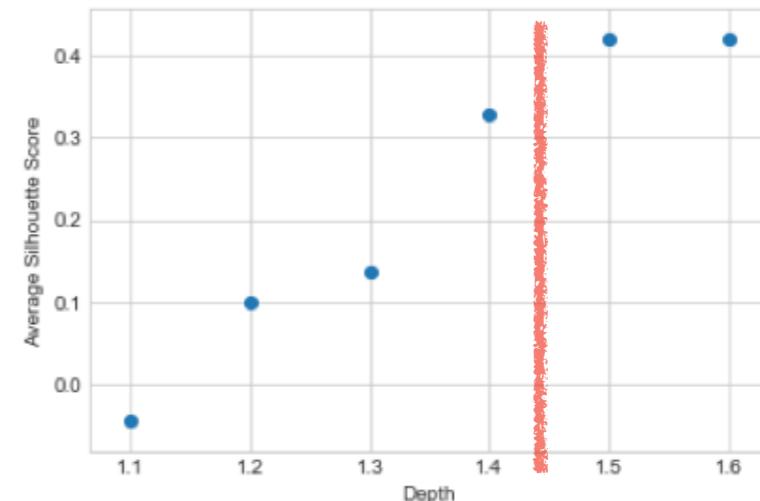
Offence Code	Description	Offence Level	Label
341	PETIT LARCENY	MISDEMEANOR	0
109	GRAND LARCENY	FELONY	0
578	HARRASSMENT 2	VIOLATION	0
344	ASSAULT 3 & RELATED OFFENSES	MISDEMEANOR	0
351	CRIMINAL MISCHIEF & RELATED OF	MISDEMEANOR	0
341	PETIT LARCENY	MISDEMEANOR	1
578	HARRASSMENT 2	VIOLATION	1
344	ASSAULT 3 & RELATED OFFENSES	MISDEMEANOR	1
351	CRIMINAL MISCHIEF & RELATED OF	MISDEMEANOR	1
235	DANGEROUS DRUGS	MISDEMEANOR	1
341	PETIT LARCENY	MISDEMEANOR	2
578	HARRASSMENT 2	VIOLATION	2
351	CRIMINAL MISCHIEF & RELATED OF	MISDEMEANOR	2
344	ASSAULT 3 & RELATED OFFENSES	MISDEMEANOR	2
109	GRAND LARCENY	FELONY	2



Hierarchical clustering was the second clustering method used that showed better silhouette scores for greater depths i.e. smaller number of clusters. Three clusters were chosen for the further study to compare the results with K-Means



Correlation - based
distance was used as
dissimilarity measure



Hierarchical clustering produced similar clusters to those in K-Means method once codes

Labels	BRONX	BROOKLYN	MANHATTAN	QUEENS	STATEN ISLAND
1	10	13	9	10	1
2	0	1	13	0	0
3	2	9	0	6	3

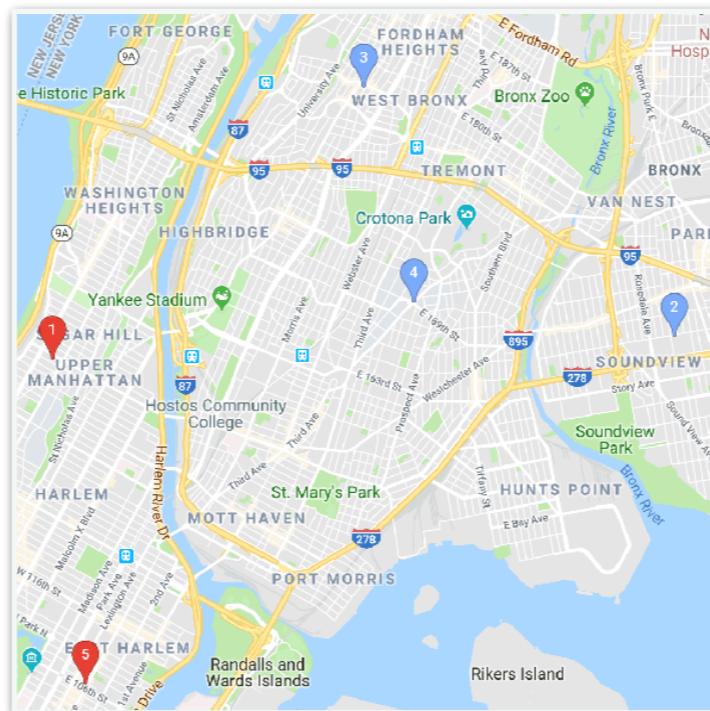
Offence Code	Description	Offence Level	Label
341	PETIT LARCENY	MISDEMEANOR	1
578	HARRASSMENT 2	VIOLATION	1
344	ASSAULT 3 & RELATED	MISDEMEANOR	1
351	CRIMINAL MISCHIEF &	MISDEMEANOR	1
235	DANGEROUS DRUGS	MISDEMEANOR	1
341	PETIT LARCENY	MISDEMEANOR	2
109	GRAND HARRASSMENT	FELONY	2
578	2	VIOLATION	2
344	ASSAULT 3 & RELATED	MISDEMEANOR	2
351	CRIMINAL MISCHIEF &	MISDEMEANOR	2
341	PETIT LARCENY	MISDEMEANOR	3
578	HARRASSMENT 2	VIOLATION	3
351	CRIMINAL MISCHIEF &	MISDEMEANOR	3
109	GRAND ASSAULT 3 &	FELONY	3
344	RELATED	MISDEMEANOR	3



A tool to study cosine similarities among precincts was created that shows on the map the location of the most similar precincts and the offence codes that contribute the most to those similarities

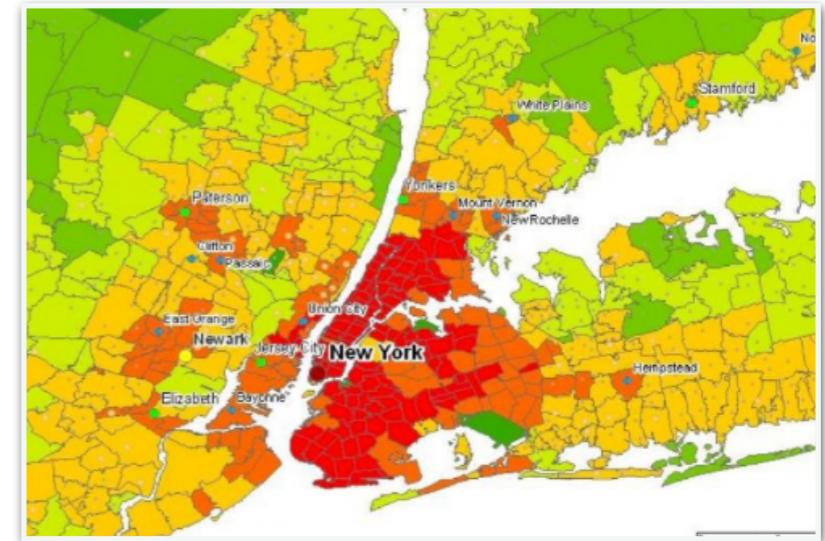
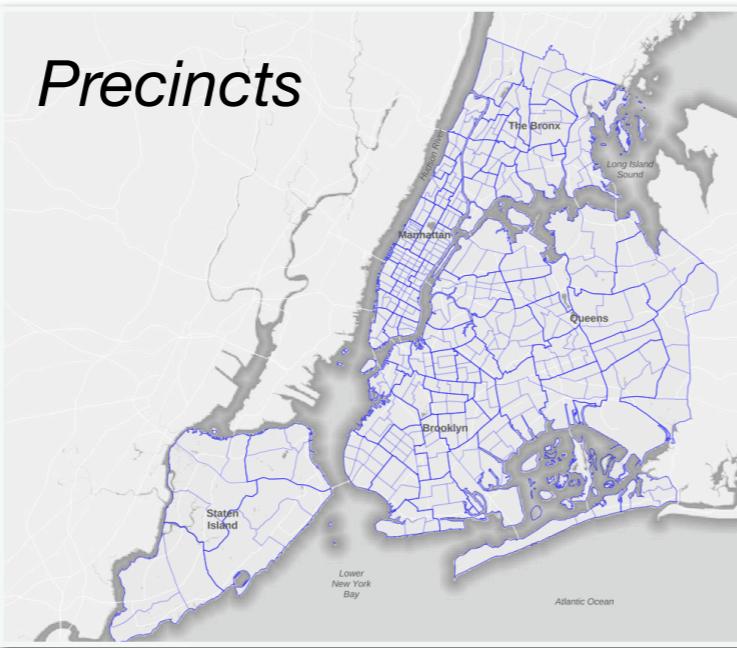
Precinct	Map order	Cosine product
30	1 (selected precinct)	1.0000
43	2	0.9923
46	3	0.9881
42	4	0.9812
23	5	0.9781

Offence Code	Description	Offence Level
235	DANGEROUS DRUGS	MISDEM EANOR
344	ASSAULT 3 & RELATED OFFENSES	MISDEM EANOR
105	ROBBERY	FELONY
351	CRIMINAL MISCHIEF & RELATED OF	MISDEM EANOR
106	FELONY ASSAULT	FELONY



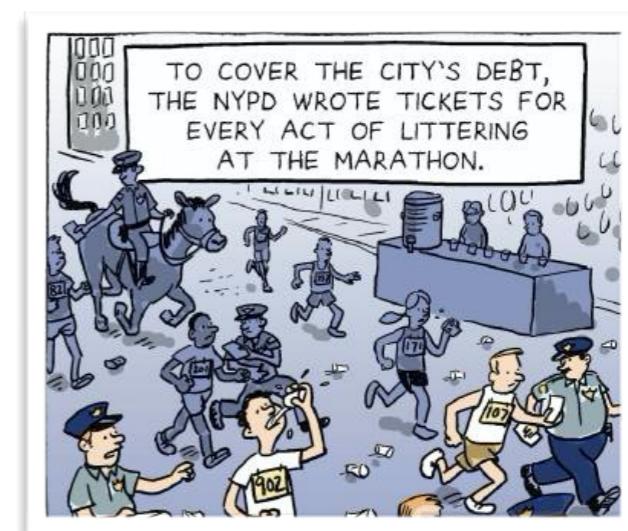
Assumptions and Limitations

- Two datasets used operate on different granularity level - precincts and neighborhoods often overlap but still differ in size and number
- Homogenous population density was assumed but it varies among precincts
- No information on number of police resources deployed in each borough or complaints dataset limitations



Different densities

Different sub-borough granularities
Neighborhood ≠ Precinct



Different police resources ?

Recommendations

- The continuation of this analysis requires more information on the data collection conditions
- Neighborhoods are more natural divisions than precincts = more useful information
- Not every type of crime is indicator of (un)safeness - selection of severe crime codes can have more practical application



Verification of
NYPD database

Neighborhood
Level Analysis

Analysis per crime type
*What crimes decrease
safety?*