Initializing System

```
1 !apt-get install openjdk-8-jdk-headless -qq > /dev/null
```

```
1 !tar xf "/content/drive/My Drive/BigDataAssignment3Files/spark-2.4.5-bin-hadoop2.7.tgz"
2 !pip install -q findspark
```

```
1 import os
2 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
3 os.environ["SPARK_HOME"] = "/content/spark-2.4.5-bin-hadoop2.7"
```

```
1 import findspark
2 findspark.init()
3 from pyspark.sql import SparkSession
4 from pyspark.context import SparkContext
5 spark = SparkSession.builder.master("local[*]").getOrCreate()
```

Reading Preprocessed Dataset

```
1 preprocessed_data = spark.read.json("hdfs://udit_gupta_1/processed_data")
```

Extract Data for EDA

```
1 preprocessed_data.columns
```

```
1 #Get top 30 crime types in complaints
2 top30_crime_type = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'C') \
4 .map(lambda row : (row['OFNS_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 complaints_crime_list = [ele[0] for ele in top30_crime_type if ele[0] is not None]
```

```
1 #Get top 30 arrests crime types
2 top30_arrests_crime_type = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'A') \
4 .map(lambda row : (row['OFNS_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 arrests_crime_list = [ele[0] for ele in top30_arrests_crime_type if ele[0] is not None]
```

```
1 #Get top 30 location types for crime complaints
2 top30_crime_locations = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'C') \
4 .map(lambda row : (row['PREM TYP DESC'],1)) \
```

```
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 complaints_location_list = [ele[0] for ele in top30_crime_locations if ele[0] is not None]
```

Generic Imports

```
1 import pandas as pd
2 import plotly.express as px
3 from pyspark.sql.functions import unix_timestamp, from_unixtime
4 from pyspark.sql import functions as F
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import folium
9 from folium.plugins import HeatMap
10 import copy
11
12 %matplotlib inline
```

## ▾ Spatial Analysis

```
1 #Filtering for Last 3 years of data on Spatial Analysis
2 arrests_pandas_df = preprocessed_data.filter(preprocessed_data['RECORD_TYPE'] == 'A').select("*",from_unixtime(unix_timestamp('ARREST_DATE', 'MM/dd/yyyy')).alias(
```

```
1 arrests_pandas_df['Latitude'] = arrests_pandas_df['Latitude'].astype(float)
2 arrests_pandas_df['Longitude'] = arrests_pandas_df['Longitude'].astype(float)
```
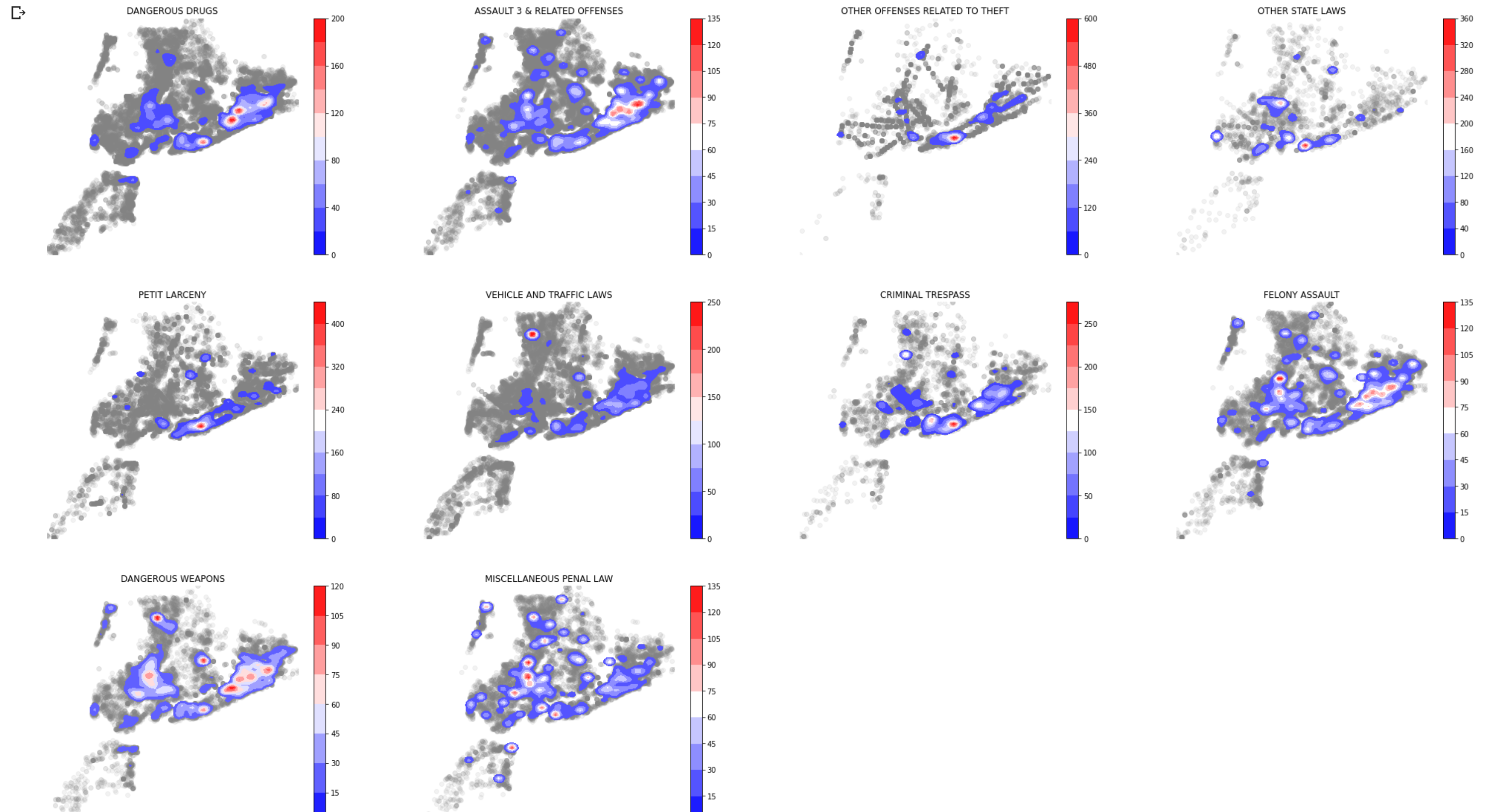
Density of top 10 Crimes in NYC

```
1 fig = plt.figure(figsize=(35,20))
2 for i, crime_type in enumerate(arrests_crime_list[0:10]):
3     ax = fig.add_subplot(int(np.ceil(float(len(arrests_crime_list[0:10])) / 4)), 4, i+1)
4     crimes_ = arrests_pandas_df[arrests_pandas_df['OFNS_DESC']==crime_type]
5
6     #Plots a scatter plot style graph as a base map
7     sns.regplot(crimes_['Latitude'], crimes_['Longitude'],
8                 fit_reg=False,
9                 scatter_kws={'alpha':.1, 'color':'grey'},
10                ax = ax)
11
12     #Plots a bivariate distribution on top of Previous Distribution
13     sns.kdeplot(crimes_['Latitude'], crimes_['Longitude'],
14                cmap="bwr",
15                bw=.005,
16                #n_levels=10,
17                cbar=True,
18                shade=True,
19                shade_lowest=False,
20                ax = ax)
```

```
20       ax = ax)
21    ax.set_title(crime_type)
22    ax.set_xlim(min(crimes_['Latitude']),max(crimes_['Latitude']))
23    ax.set_ylim(min(crimes_['Longitude']),max(crimes_['Longitude']))
24    ax.set_axis_off()
25 plt.show()
```



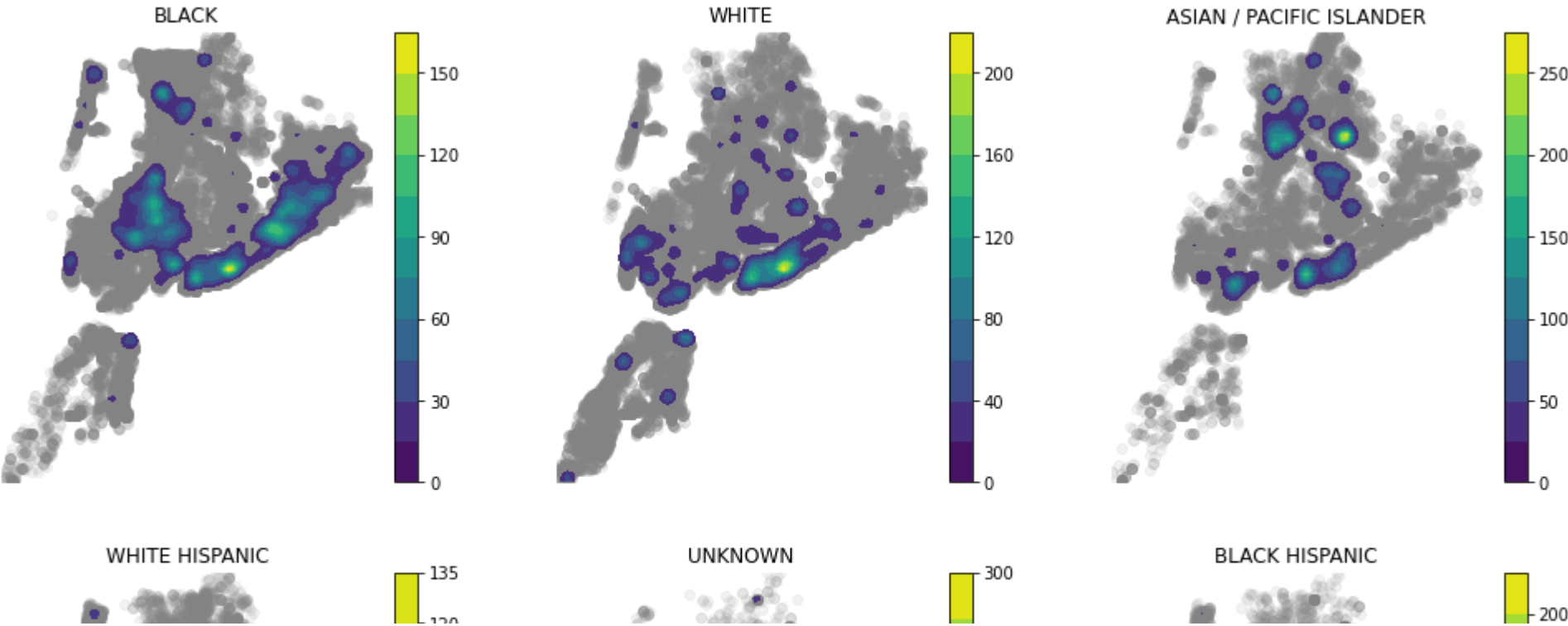Density Plot for Arrests of Different Ethnicity across NYC

```
1 fig = plt.figure(figsize=(17,17))
2 for i, race_type in enumerate(arrests_pandas_df.SUSP_RACE.unique()):
```

```
 3        ax = fig.add_subplot(int(np.ceil(float(len(arrests_pandas_df.SUSP_RACE.unique())) / 3)), 3, i+1)
 4        race_ = arrests_pandas_df[arrests_pandas_df.SUSP_RACE == race_type]
 5
 6        #Plots a scatter plot style graph as a base map
 7        sns.regplot(race_['Latitude'], race_['Longitude'],
 8                    fit_reg=False,
 9                    scatter_kws={'alpha':.1, 'color':'grey'},
10                    ax = ax)
11
12        #Plots a bivariate distribution on top of Previous Distribution
13        sns.kdeplot(race_['Latitude'], race_['Longitude'],
14                     cmap="viridis",
15                     bw=.005,
16                     #n_levels=10,
17                     cbar=True,
18                     shade=True,
19                     shade_lowest=False,
20                     ax = ax)
21    ax.set_title(race_type)
22    ax.set_xlim(min(race_['Latitude']),max(race_['Latitude']))
23    ax.set_ylim(min(race_['Longitude']),max(race_['Longitude']))
24    ax.set_axis_off()
25 plt.show()
```
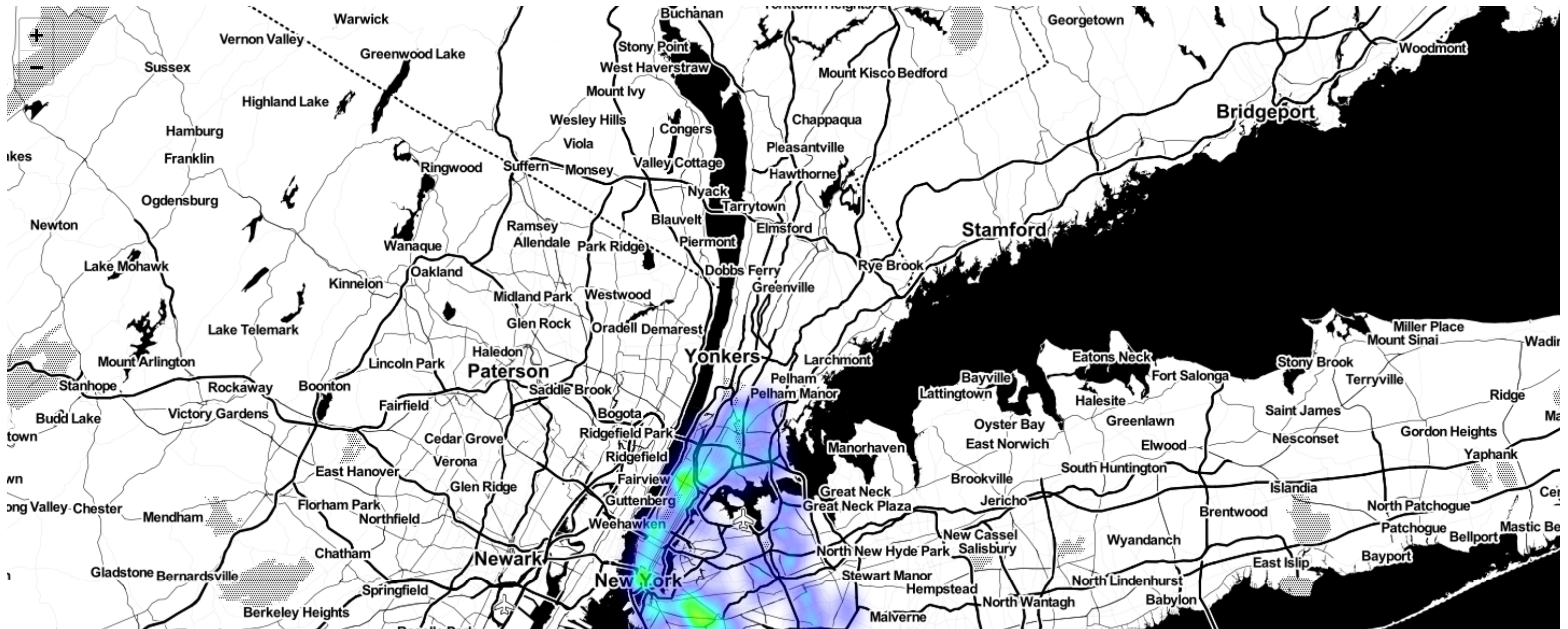
BLACK WHITE ASIAN / PACIFIC ISLANDER

WHITE HISPANIC UNKNOWN BLACK HISPANIC

▾ Heat Map for A Closer look at some Violent Crimes in NYC

Harrassment

```
1 # Harrassment crime map
2 crime_map = folium.Map(location=[40.93, -73.91],
3                        tiles = "Stamen Toner",
4                        zoom_start = 10)
5 data_heatmap = []
6 # Add data for heatmap
7 for index, row in arrests_pandas_df[arrests_pandas_df['OFNS_DESC'] == 'HARRASSMENT 2'].iterrows():
8   if(row['Latitude'] is None or row['Longitude'] is None):
9       continue
10  data_heatmap.append([float(row['Latitude']), float(row['Longitude'])])
11
12 HeatMap(data_heatmap, radius=10).add_to(crime_map)
13 crime_map
```
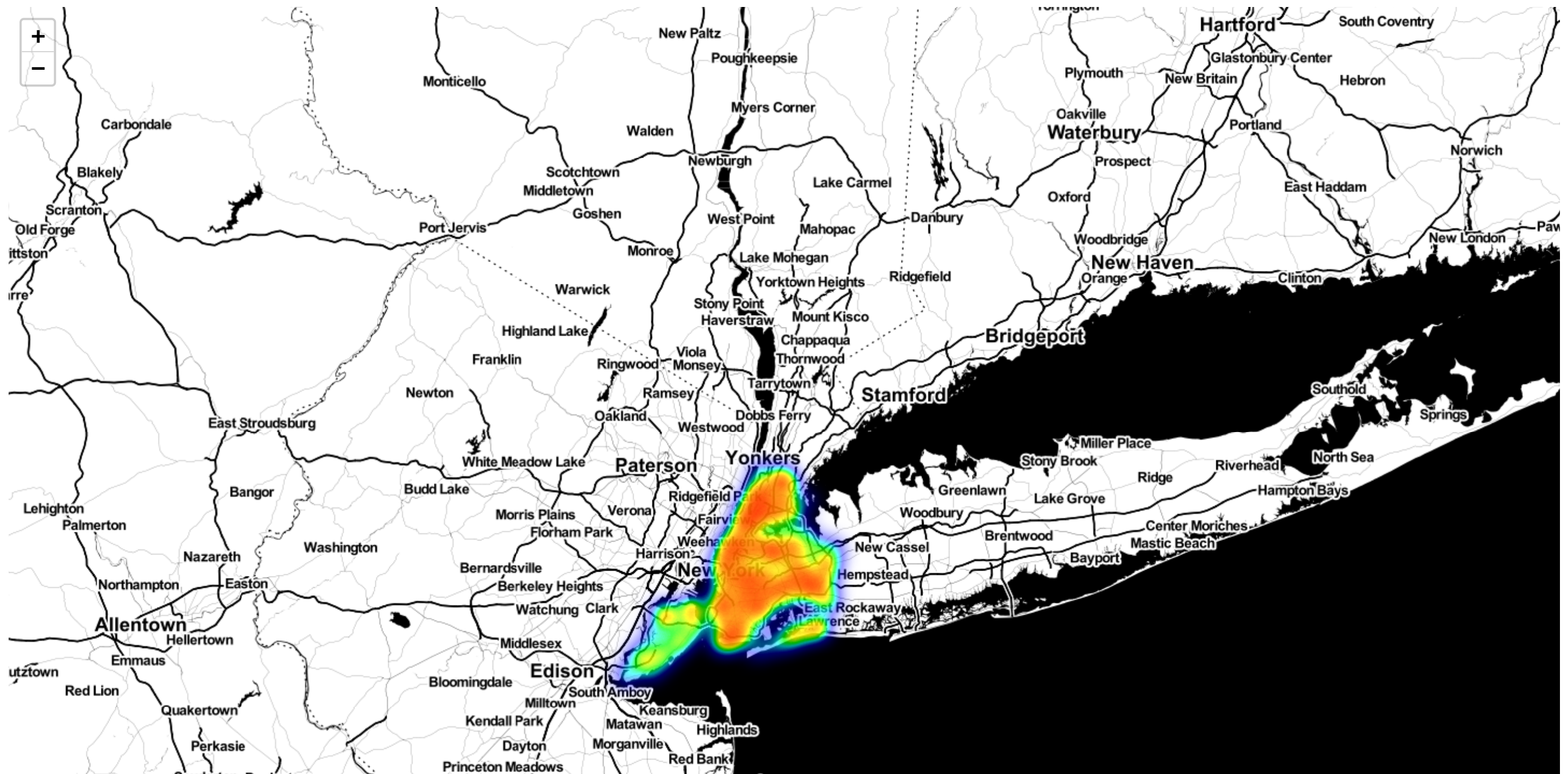
↪

Assault

```
1 # Harrassment crime map
2 crime_map = folium.Map(location=[40.93, -73.91],
3                        tiles = "Stamen Toner",
4                   zoom_start = 10)
5 data_heatmap = []
6 # Add data for heatmap
7 for index, row in arrests_pandas_df[arrests_pandas_df['OFNS_DESC'] == 'ASSAULT 3 & RELATED OFFENSES'].iterrows():
8   if(row['Latitude'] is None or row['Longitude'] is None):
9       continue
10  data_heatmap.append([float(row['Latitude']), float(row['Longitude'])])
11
12 HeatMap(data_heatmap, radius=10).add_to(crime_map)
13 crime_map
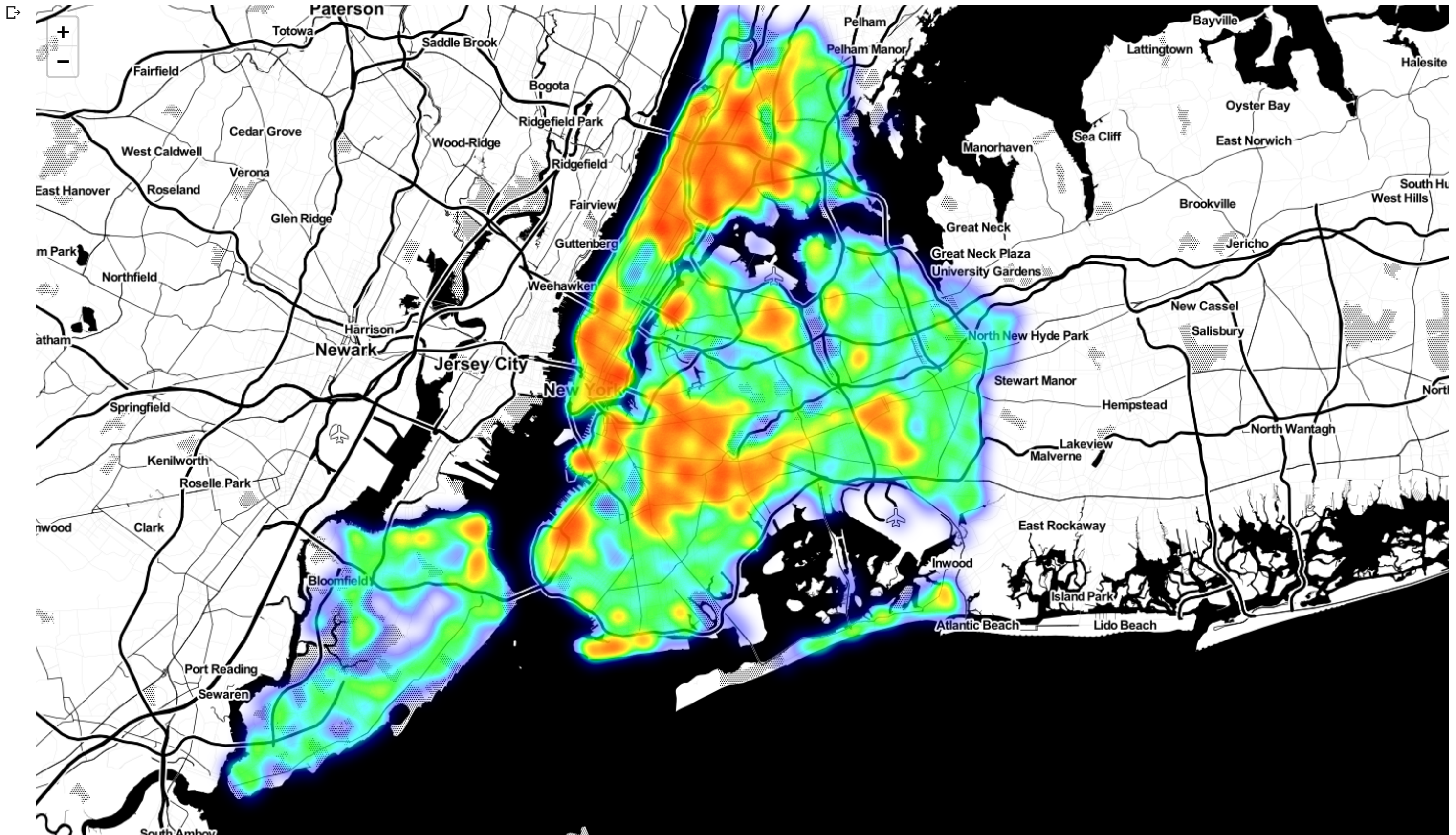```

⤷

Dangerous Drugs

```
1  # Harrassment crime map
2  crime_map = folium.Map(location=[40.7, -73.91],
3                         tiles = "Stamen Toner",
4                         zoom_start = 10)
5  data_heatmap = []
6  # Add data for heatmap
7  for index, row in arrests_pandas_df[arrests_pandas_df['OFNS_DESC'] == 'DANGEROUS DRUGS'].iterrows():
8    if(row['Latitude'] is None or row['Longitude'] is None):
9        continue
10   data_heatmap.append([float(row['Latitude']), float(row['Longitude'])])
11
12 HeatMap(data_heatmap, radius=10).add_to(crime_map)
13 crime_map
```

Stratified Sampling for Large Data

```
1  # Stratified Sampling
2  race_fraction = {
3      "WHITE":0.01,
4      "BLACK":0.01,
5      "BLACK HISPANIC":0.01,
6      "WHITE HISPANIC":0.01,
7      "UNKNOWN":0.
```

```
 7        UNKNOWN :0,
 8        "OTHER":0,
 9        "AMERICAN INDIAN/ALASKAN NATIVE":0.01,
10        "ASIAN / PACIFIC ISLANDER":0.01
11 }
12
13 sampled_data = preprocessed_data \
14 .filter(preprocessed_data['RECORD_TYPE'] == 'A') \
15 .sampleBy("SUSP_RACE", fractions=race_fraction, seed=0).toPandas()
16
17 sampled_data = sampled_data[["Latitude", "Longitude", "SUSP_RACE"]]\
18                 .astype({'Latitude': 'double', 'Longitude': 'double'})
```

Ethinicity Distribution of Arrests across NYC

```
1 map_box_token = "_Security_Key_Here_"
2 px.set_mapbox_access_token(map_box_token)
3 df = px.data.carshare()
4 fig = px.scatter_mapbox(sampled_data, lat="Latitude", lon="Longitude", color="SUSP_RACE",
5                 color_continuous_scale=px.colors.cyclical.IceFire, zoom=11 ,width=1200, height=768)
6 fig.show()
```

⤷