NYU

# The Anatomy of NYPD Crime Records:

## Race, Gender, Politics and More

By: Jiaming Chen, Zizhen Chen, Jiayuan Huang, Yu Wu

# Goals

- Intergroup Solidarity, Felt and Shared Foreignness between Asians and Hispanics in the U.S ✅
- They are under similar social pressure  **?**
  Similar crime behavior patterns in NYC  (correlation based on crime amounts and types distribution ) ❌

- **No hypothesis.**
- **All Races**
- **Various Analysis**
- **Direct Explaining**

# Data

Unbiased data  = How?

- Not arrest data but Complaint Data
- Omit minor crimes such as marijuwana, loitering, minor valued theft

Cleaned dataset = what columns?

- Race, Gender, Age, Crime type, Victim, Location…..

NYU

# Methodologies
# &
# Corresponding Results

# Clustering Analysis

Smallest Cluster for races

| day | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 | 6.0 | labelk |
|---|---|---|---|---|---|---|---|---|
| OFNS_DESC | | | | | | | | |
| HARRASSMENT 2 | 1576.0 | 1613.0 | 1578.0 | 1589.0 | 1604.0 | 1676.0 | 1638.0 | 1 |

Hispanic

| day | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 | 6.0 | labelk |
|---|---|---|---|---|---|---|---|---|
| OFNS_DESC | | | | | | | | |
| HARRASSMENT 2 | 397.0 | 377.0 | 370.0 | 385.0 | 422.0 | 391.0 | 424.0 | 1 |

Asian

| day | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 | 6.0 | labelk |
|---|---|---|---|---|---|---|---|---|
| OFNS_DESC | | | | | | | | |
| HARRASSMENT 2 | 934.0 | 906.0 | 970.0 | 901.0 | 905.0 | 880.0 | 851.0 | 1 |

White

| day | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 | 6.0 | labelk |
|---|---|---|---|---|---|---|---|---|
| OFNS_DESC | | | | | | | | |
| HARRASSMENT 2 | 2782.0 | 2694.0 | 2759.0 | 2802.0 | 2761.0 | 2667.0 | 2637.0 | 1 |

Black

Clustering using K-means with silhouette score

```
range_n_clusters = [2, 3, 4, 5, 6]

for n_clusters in range_n_clusters:

    clusterer = KMeans(n_clusters=n_clusters, random_state=0).fit(table)

    silhouette_avg = silhouette_score(table, clusterer.labels_)

    print(
        "For n_clusters =",
        n_clusters,
        "The average silhouette_score is :",
        silhouette_avg,
    )
```
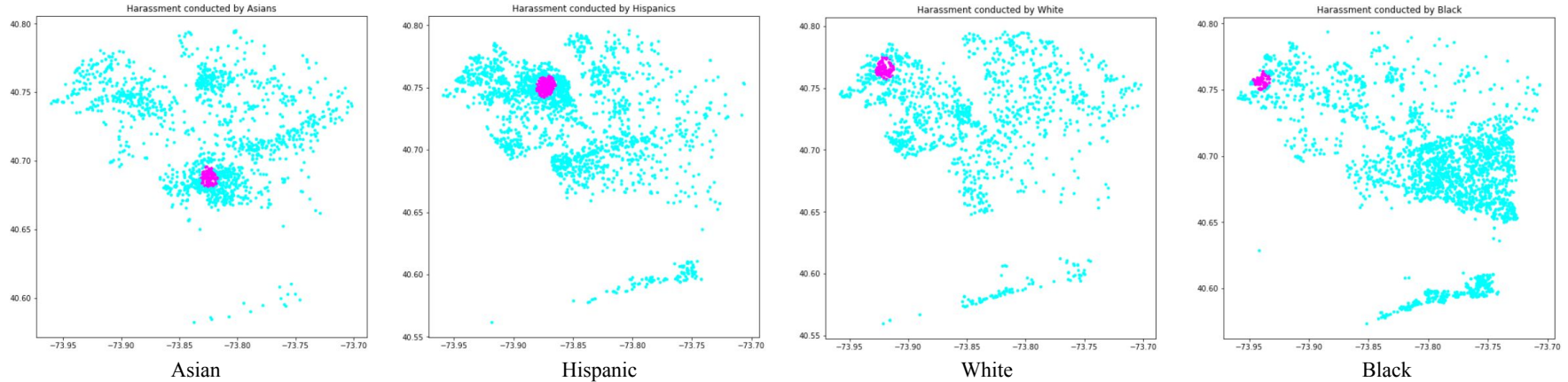
```
For n_clusters = 2 The average silhouette_score is : 0.9109035883498503
For n_clusters = 3 The average silhouette_score is : 0.789415217948282
For n_clusters = 4 The average silhouette_score is : 0.7210837816951361
For n_clusters = 5 The average silhouette_score is : 0.7092426089710087
For n_clusters = 6 The average silhouette_score is : 0.6346130603227677
```
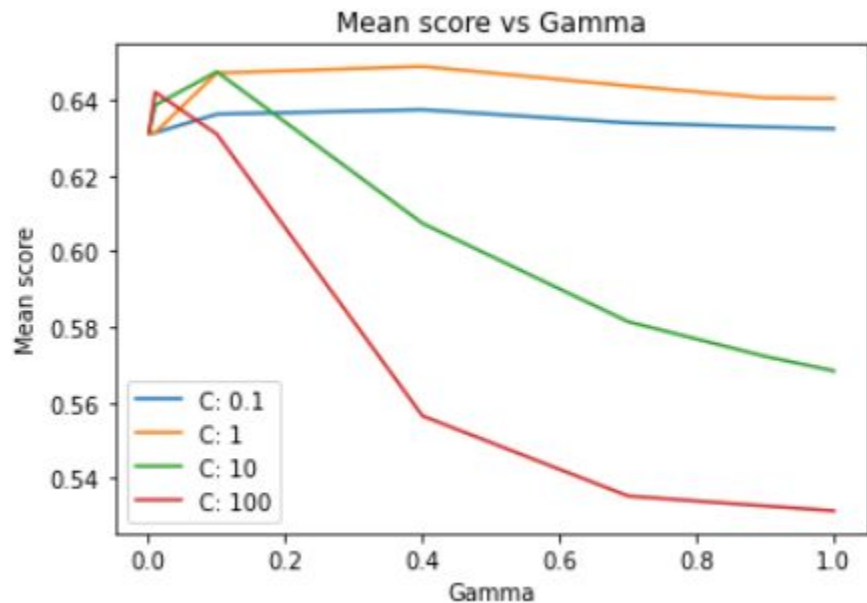
- As we can see, we will select 2 as the number of clusters
- The model cluster the harassment 2 into the smallest cluster for all the races

NYU

# DBScan Analysis



| Asian | Hispanic | White | Black |

- The lower part of Queens has the most serious asian harassment problem.
- Each race's most serious harassment problem happened in different places in Queens.

# SVM

## Mean score vs Gamma



```
grid = GridSearchCV(svm.SVC(), param_grid=param_grid, cv=2)

grid.fit(X_train, Y_train)

print("The best classifier is: ", grid.best_estimator_)
```

The best classifier is:  SVC(C=1, gamma=0.4)

```
# plot the scores of the grid
# grid_scores_ contains parameter settings and scores
ypred1 = grid.predict(X_test)
print("Out of sample, rbf svm successfully predicts {} percent of the data".format(accuracy_score(Y_test,ypred1)
```
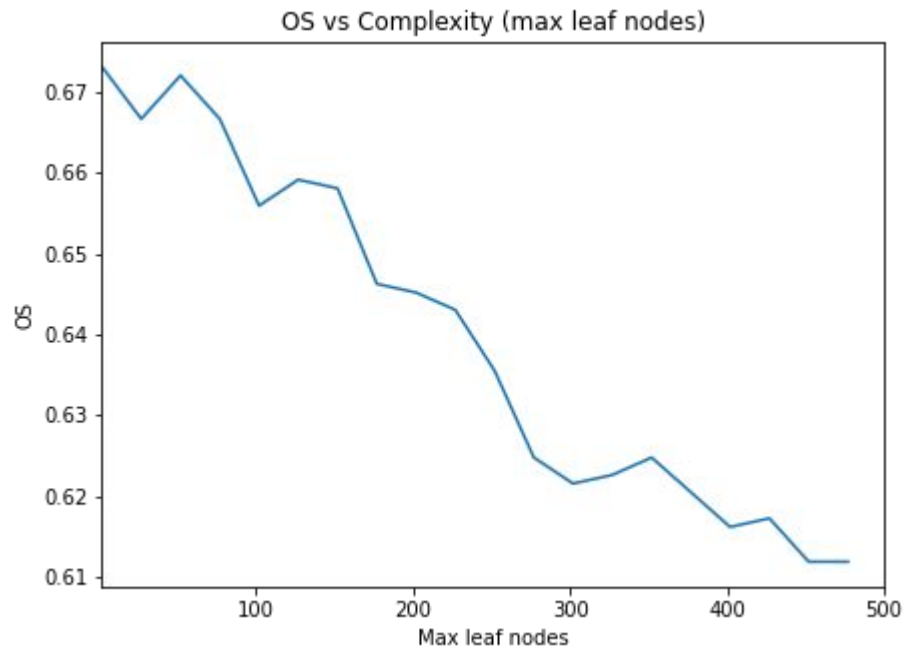
Out of sample, rbf svm successfully predicts 0.6562842528315674 percent of the data

```
grid.cv_results_['mean_test_score']
```

```
array([0.63242902,  0.63283395,  0.63395874,  0.63742294,  0.63625324,
       0.63134926,  0.63134926,  0.64039226,  0.64057224,  0.64372151,
       0.64885036,  0.64709588,  0.63134926,  0.63134926,  0.56836296,
       0.57218715,  0.58136519,  0.60736948,  0.64745568,  0.63863774,
       0.63134926,  0.53133583,  0.53264064,  0.53524996,  0.55644048,
       0.63094431,  0.64205699,  0.63130427])
```

# Decision Tree

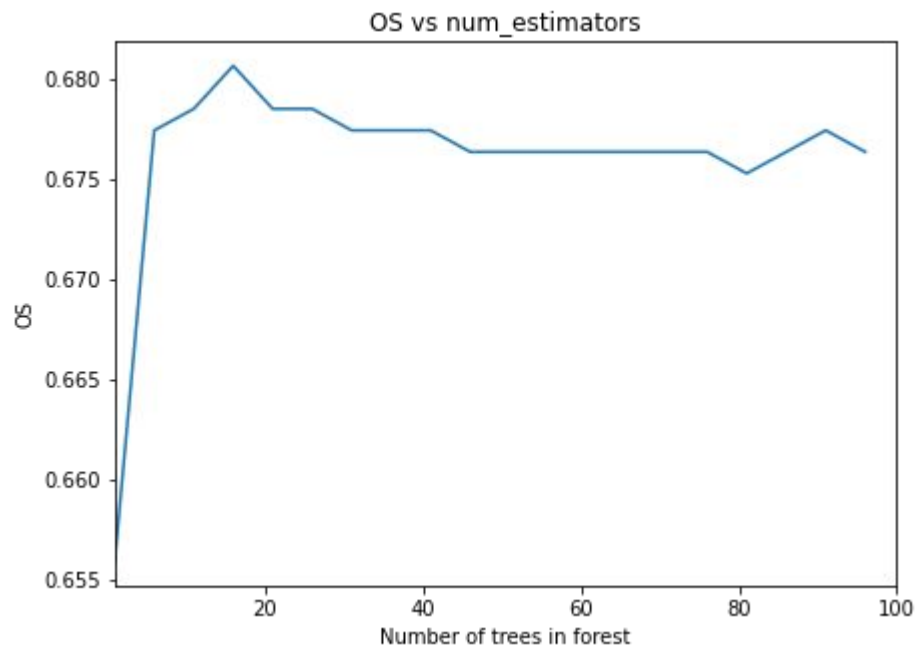OS vs Complexity (max leaf nodes)



```
# your code here
from sklearn.model_selection import GridSearchCV

param_grid = {'max_leaf_nodes':range(2,500,25)}
dt=DecisionTreeClassifier(random_state=42)
gr=GridSearchCV(dt,param_grid=param_grid,scoring='accuracy')
rs=gr.fit(X_train,y_train)
print(rs.best_params_)
print(rs.score(X_test,y_test))
```
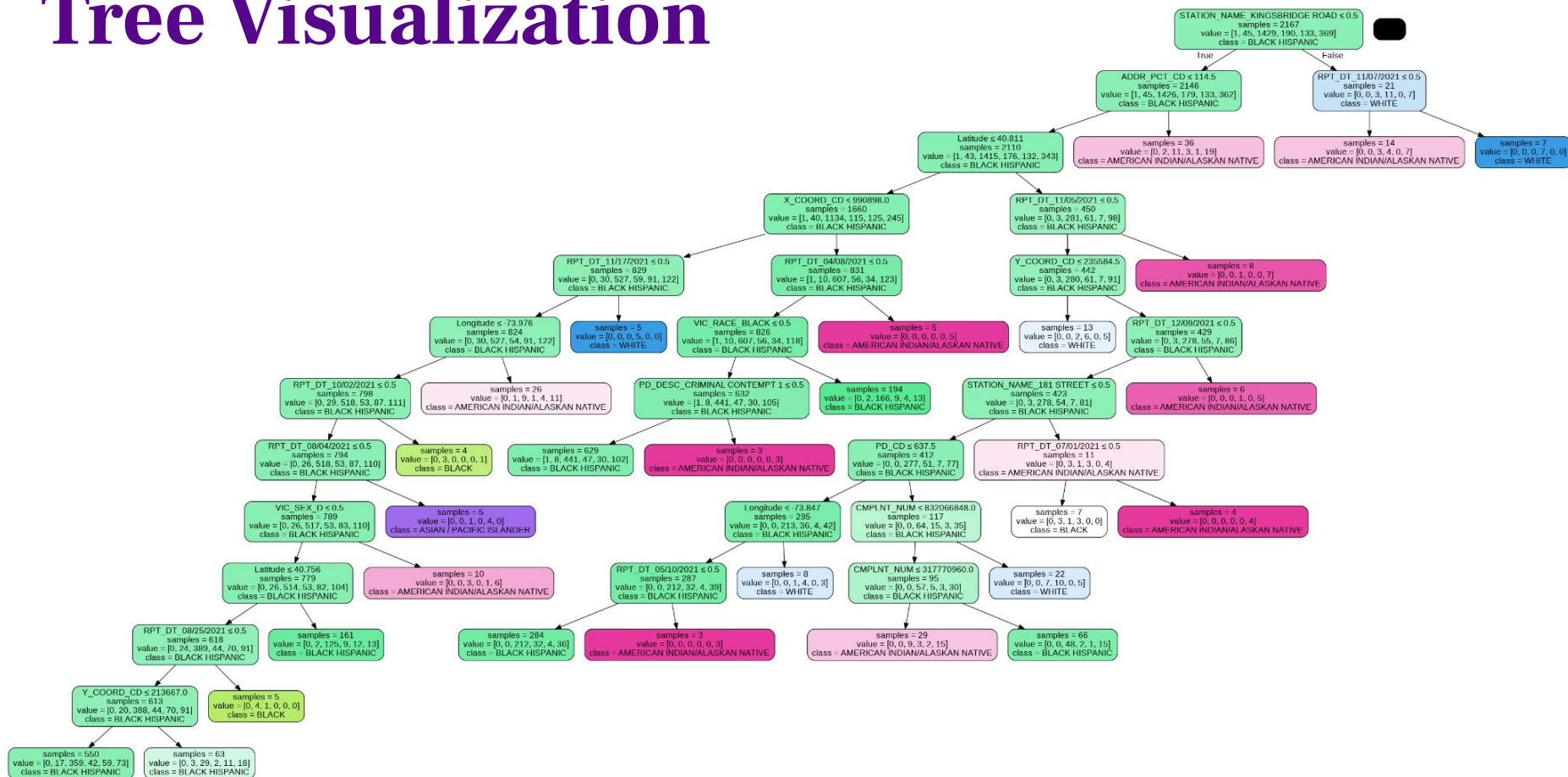
```
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection
  UserWarning,
{'max_leaf_nodes': 27}
0.6666666666666666
```

NYU

# Random Forest



OS vs num_estimators

```
param_grid = {'n_estimators':range(1,100,5)}
rf = RandomForestClassifier(n_jobs=-1,max_leaf_nodes=27, random_state = 42)
gs = GridSearchCV(rf,param_grid=param_grid,scoring='accuracy')
rs = gs.fit(X_train,y_train)
pred=rs.predict_proba(X_test)[:,1]
print(rs.best_params_)
print(rs.score(X_test,y_test))
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_split.py:680
  UserWarning,
{'n_estimators': 11}
0.678494623655914
```

# Tree Visualization

# Anomaly Detection

**Result of Gaussian Mixture:**
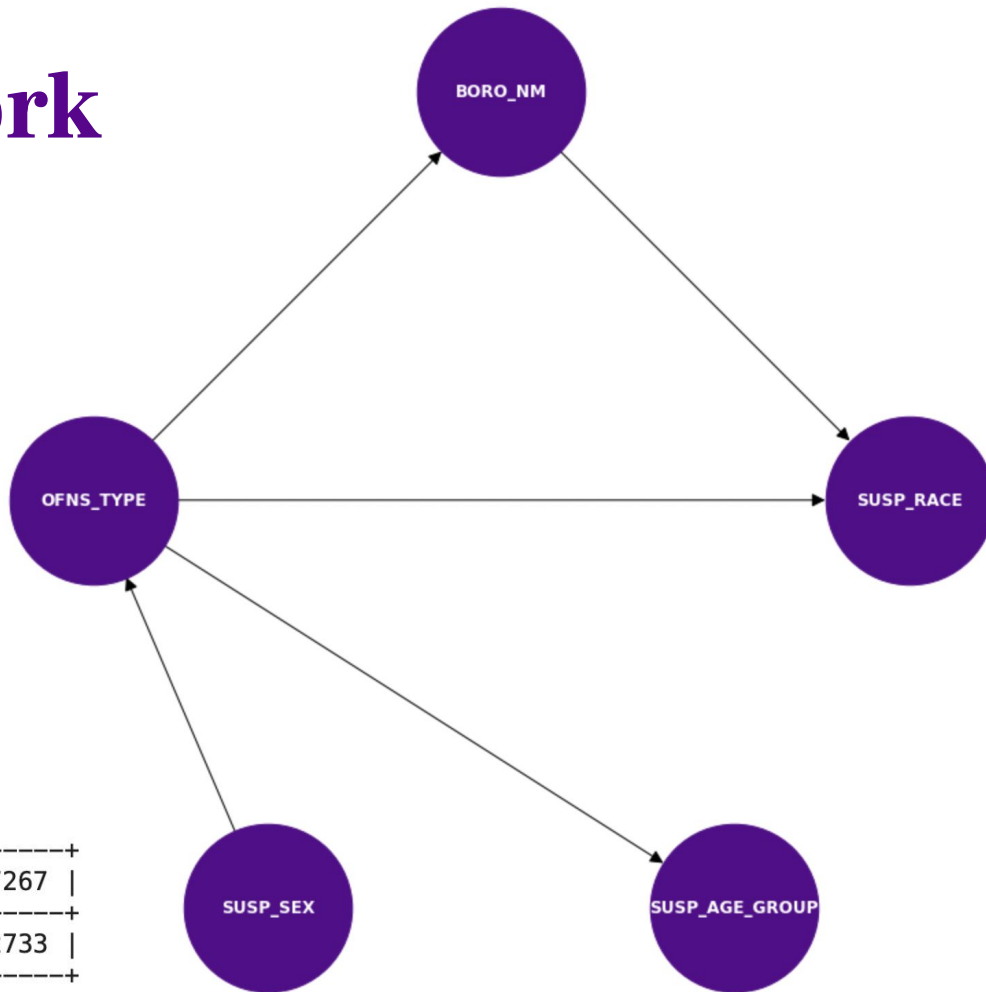
```
                                        score
OFNS_DESC
PROSTITUTION & RELATED OFFENSES  -6.653364
UNLAWFUL POSS. WEAP. ON SCHOOL   -4.762023
CHILD ABANDONMENT/NON SUPPORT    -4.579684
INTOXICATED/IMPAIRED DRIVING     -3.662789
FELONY SEX CRIMES                -2.970356
```

**Result of K-Means:**

```
                                      distance
OFNS_DESC
PROSTITUTION & RELATED OFFENSES   1.480624
INTOXICATED/IMPAIRED DRIVING      1.065731
UNLAWFUL POSS. WEAP. ON SCHOOL    1.054180
CHILD ABANDONMENT/NON SUPPORT     1.013794
HOMICIDE-NEGLIGENT-VEHICLE        0.890717
```

NYU

# Bayesian Network



```
+-----------------------------------+--------------+--------------+
| SUSP_SEX                          | SUSP_SEX(F)  | SUSP_SEX(M)  |
+-----------------------------------+--------------+--------------+
| OFNS_DESC(FELONY ASSAULT)         | 0.21         | 0.22         |
+-----------------------------------+--------------+--------------+
| OFNS_DESC(HARRASSMENT 2)          | 0.71         | 0.6          |
+-----------------------------------+--------------+--------------+
| OFNS_DESC(MISCELLANEOUS PENAL LAW)| 0.08         | 0.17         |
+-----------------------------------+--------------+--------------+
```

```
CPD of SUSP_SEX:
+--------------+----------+
| SUSP_SEX(F)  | 0.27267  |
+--------------+----------+
| SUSP_SEX(M)  | 0.72733  |
+--------------+----------+
```

NYU

# Conclusions

- Victims: Asians, American Indians are similar, Hispanics and Blacks are similar when they are victims

- SVM and Random Forest: decent accuracy could be improved by tuning more precisely

- Servere harassment is most popular for all races. Crime locations are discrete for all races

- Offense types affect community's formation