

Understanding and Predicting Crime in San Francisco

Capstone Project

Elton Yeo, DSI 13, General Assembly

The City of San Francisco and San Francisco Police Dept have limited resources...

- Where should they focus their efforts to reduce crime in the city and increase safety?
- How can they modify their policy or operational approaches to reduce crime?



This project has two objectives:

1. **Develop insights into the areas with the most crimes**, and consider alternative policy and operational measures to reduce crime



2. **Predict the number of preventable crime** given a specific zipcode, day of the week, and hour, which will help the police plan their patrol schedules and resources to reduce crime



Data



SFGovCoordinator's PortalAboutHelp

DataSF

OPEN DATASHOWCASEPUBLISHINGACADEMYRESOURCESBLOG

ExploreBrowse DataOpen Data StatsDevelopers



Sign In

Police Department Incident Reports: 2018 to Present

Public Safety

View DataVisualize ▾ExportAPI...

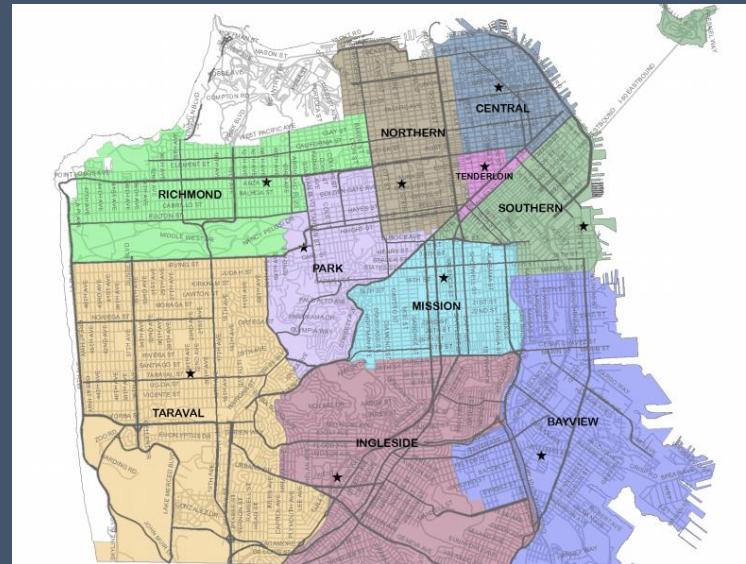
For a detailed overview, see: <https://support.datasf.org/help/police-department-incident-reports-2018-to-present-overview>

Updated
April 22, 2020

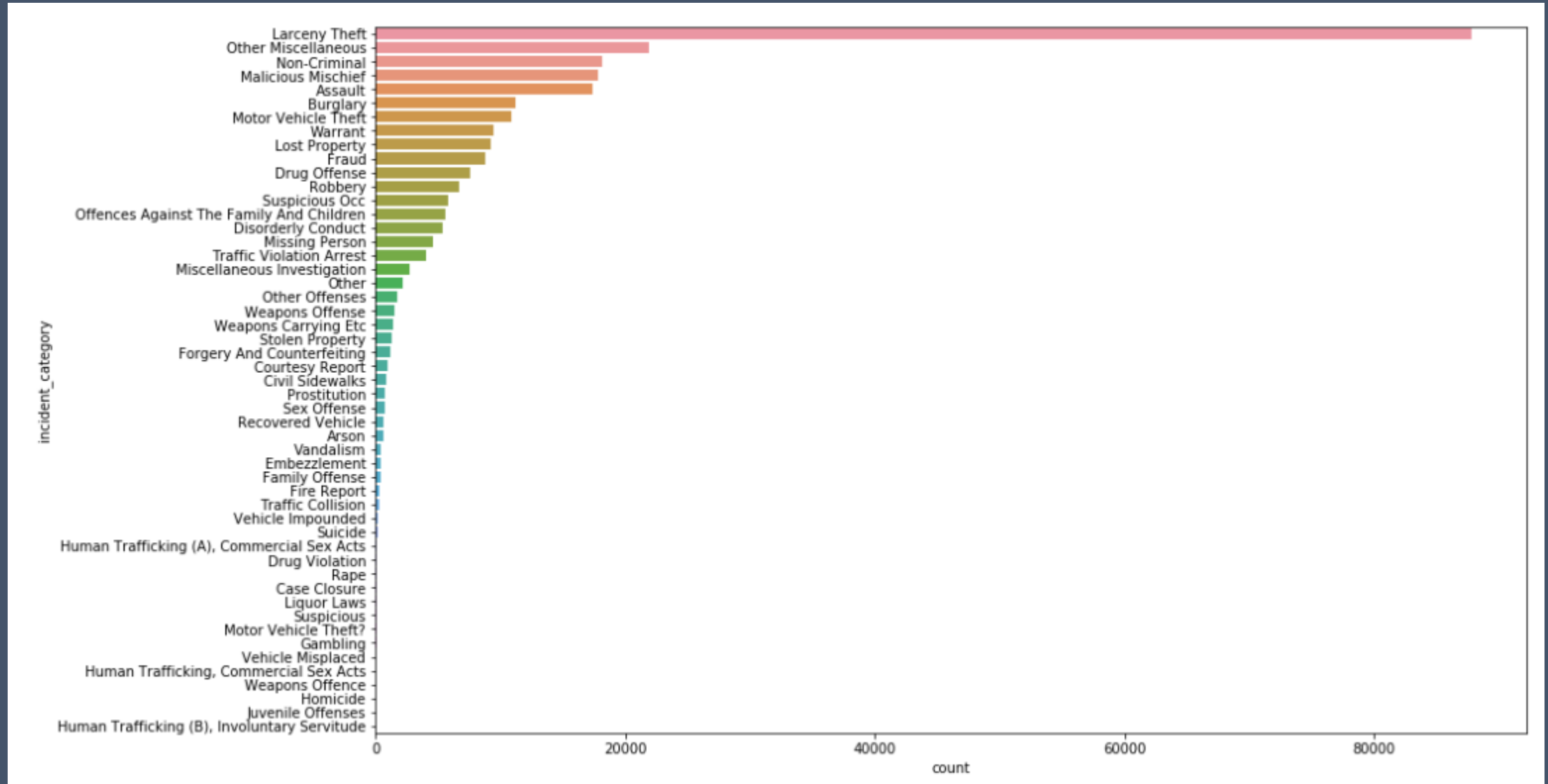
This dataset includes police incident reports filed by officers and by individuals through self-service online reporting for non-emergency cases. Reports included are those for incidents that occurred starting January 1, 2018 onward and have been approved by a supervising officer.

Making sense of our data e.g. locations

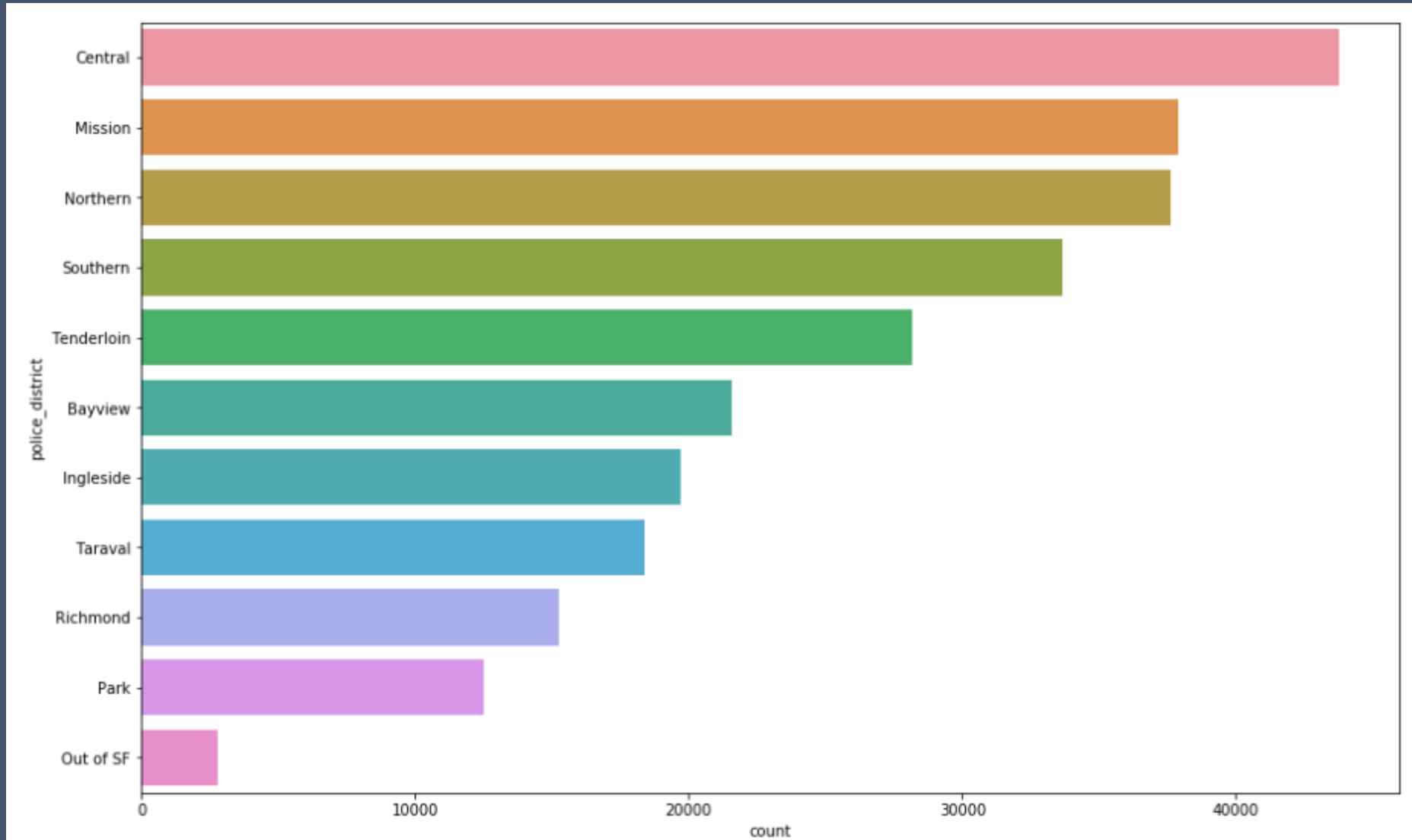
- Latitude, longitude, intersection, sf find neighbourhood, analysis neighbourhood, analysis neighbourhoods, current police district, current supervisor districts, supervisor district – what do these mean and are they important for our analysis and prediction?



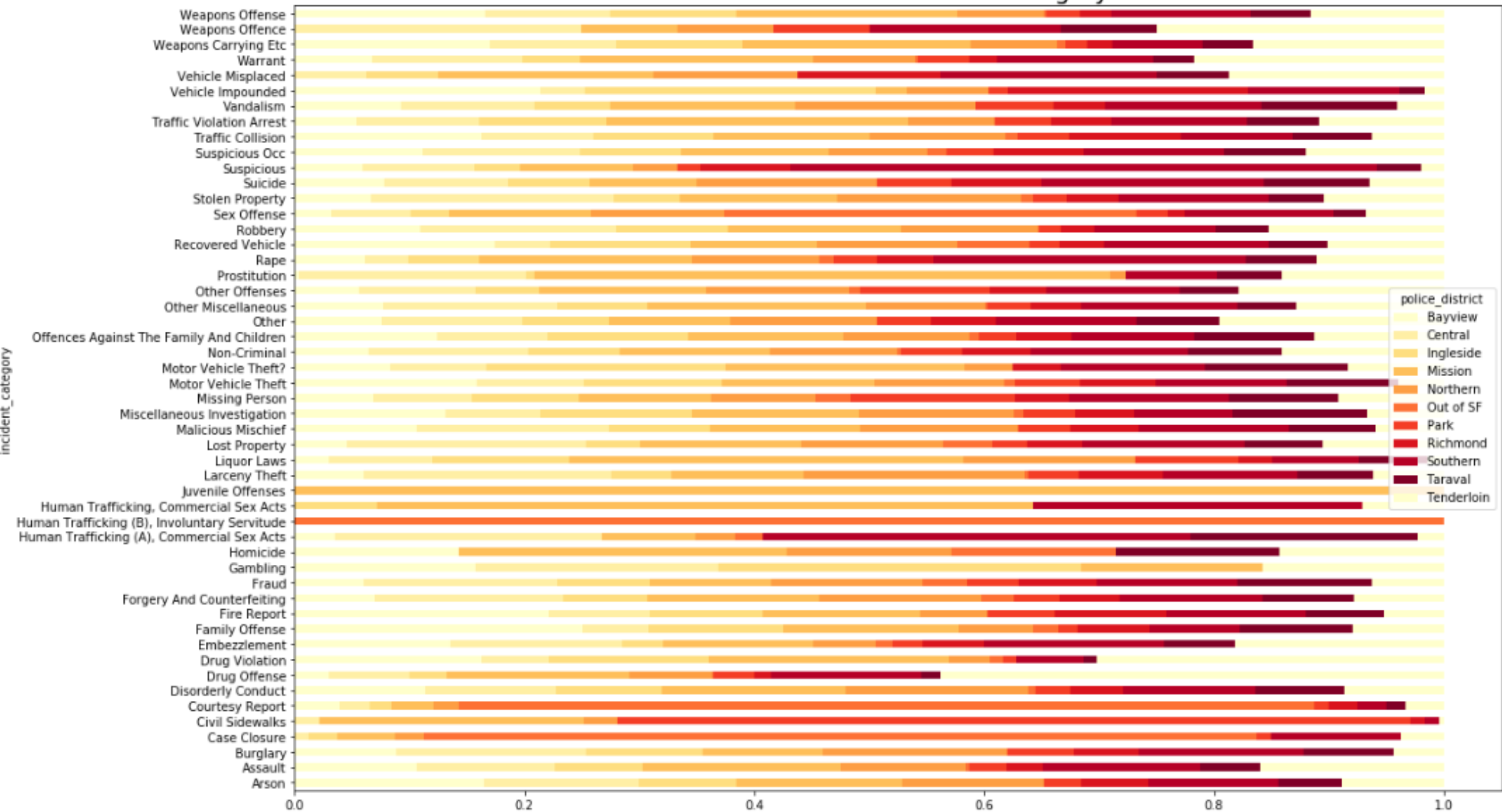
Theft occurred the most frequently – but how do we treat these different categories of crime?



Central police district had the highest number of crimes



Police District vs Incident Category



Can we group crime into more meaningful categories for analysis and prediction?

Preventable Crime

- 'Larceny Theft', 'Malicious Mischief', 'Burglary', 'Motor Vehicle Theft', 'Robbery', 'Stolen Property', 'Disorderly Conduct', 'Vandalism', 'Arson', 'Motor Vehicle Theft', 'Suspicious Occ', 'Offences Against The Family And Children', 'Suspicious', 'Family Offense', 'Prostitution', 'Civil Sidewalks'

Violent Crime

- 'Assault', 'Weapons Offense', 'Weapons Carrying Etc', 'Rape', 'Weapons Offence', 'Homicide'

Drug Crime

- 'Drug Offense', 'Drug Violation'

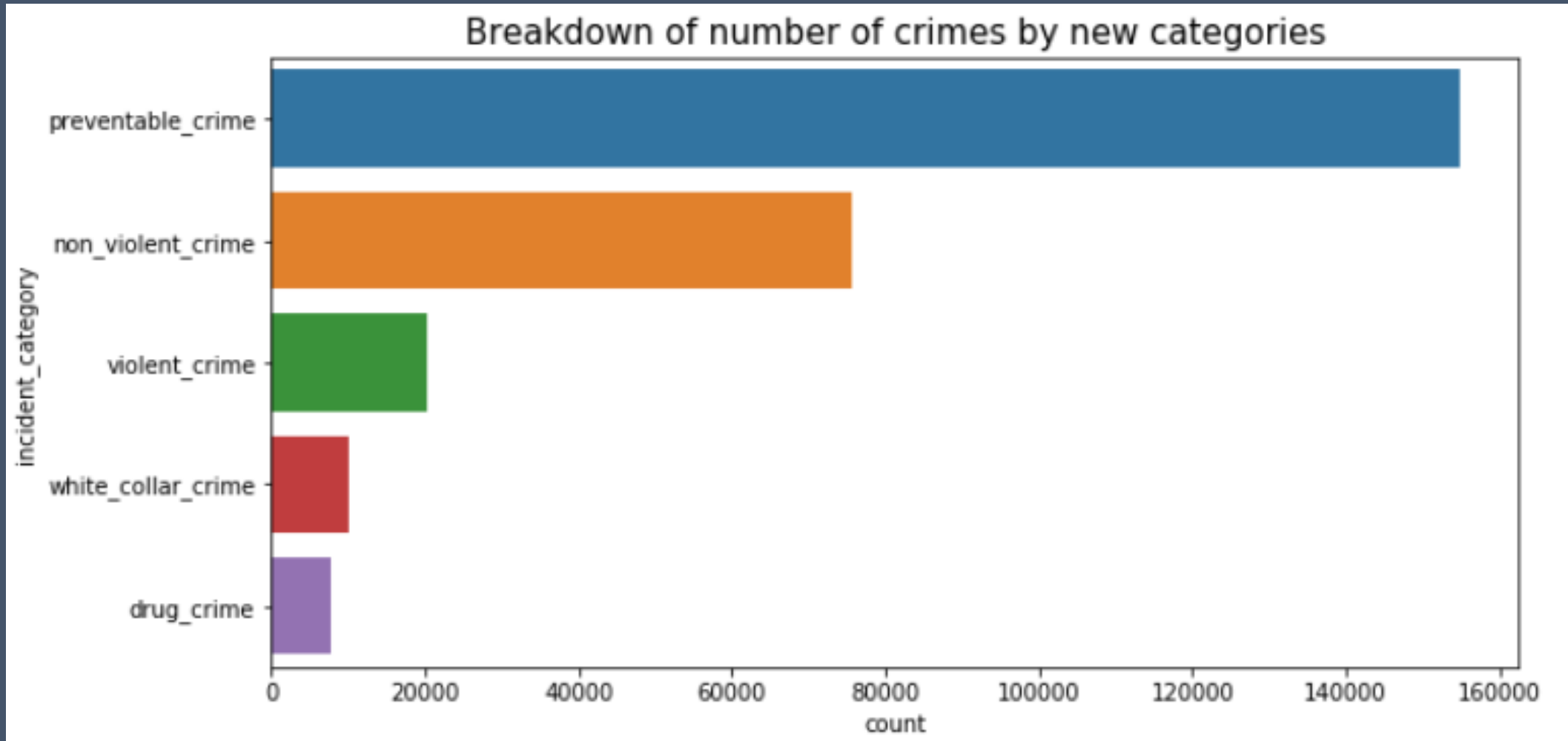
White-Collar Crime

- 'Fraud', 'Forgery And Counterfeiting', 'Embezzlement'

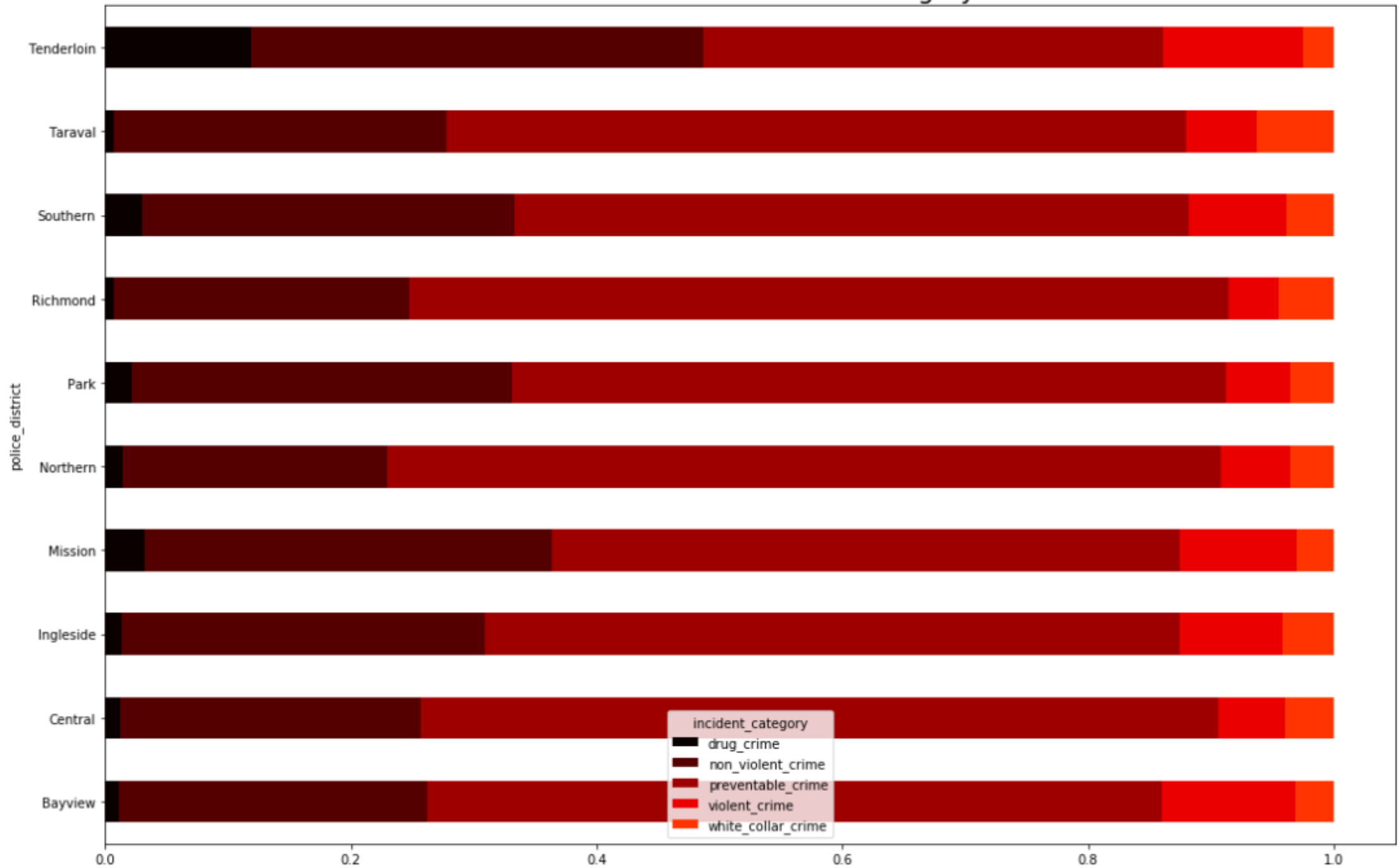
Non-Violent Crime

- 'Lost Property', 'Missing Person', 'Recovered Vehicle', 'Traffic Collision', 'Vehicle Impounded', 'Suicide', 'Vehicle Misplaced', 'Traffic Violation Arrest', 'Non-Criminal', 'Other Miscellaneous', 'Warrant', 'Gambling', 'Human Trafficking, Commercial Sex Acts', 'Human Trafficking (B), Involuntary Servitude', 'Juvenile Offenses', 'Liquor Laws', 'Case Closure', 'Human Trafficking (A), Commercial Sex Acts', 'Fire Report', 'Sex Offense', 'Courtesy Report', 'Other Offenses', 'Other', 'Miscellaneous Investigation'

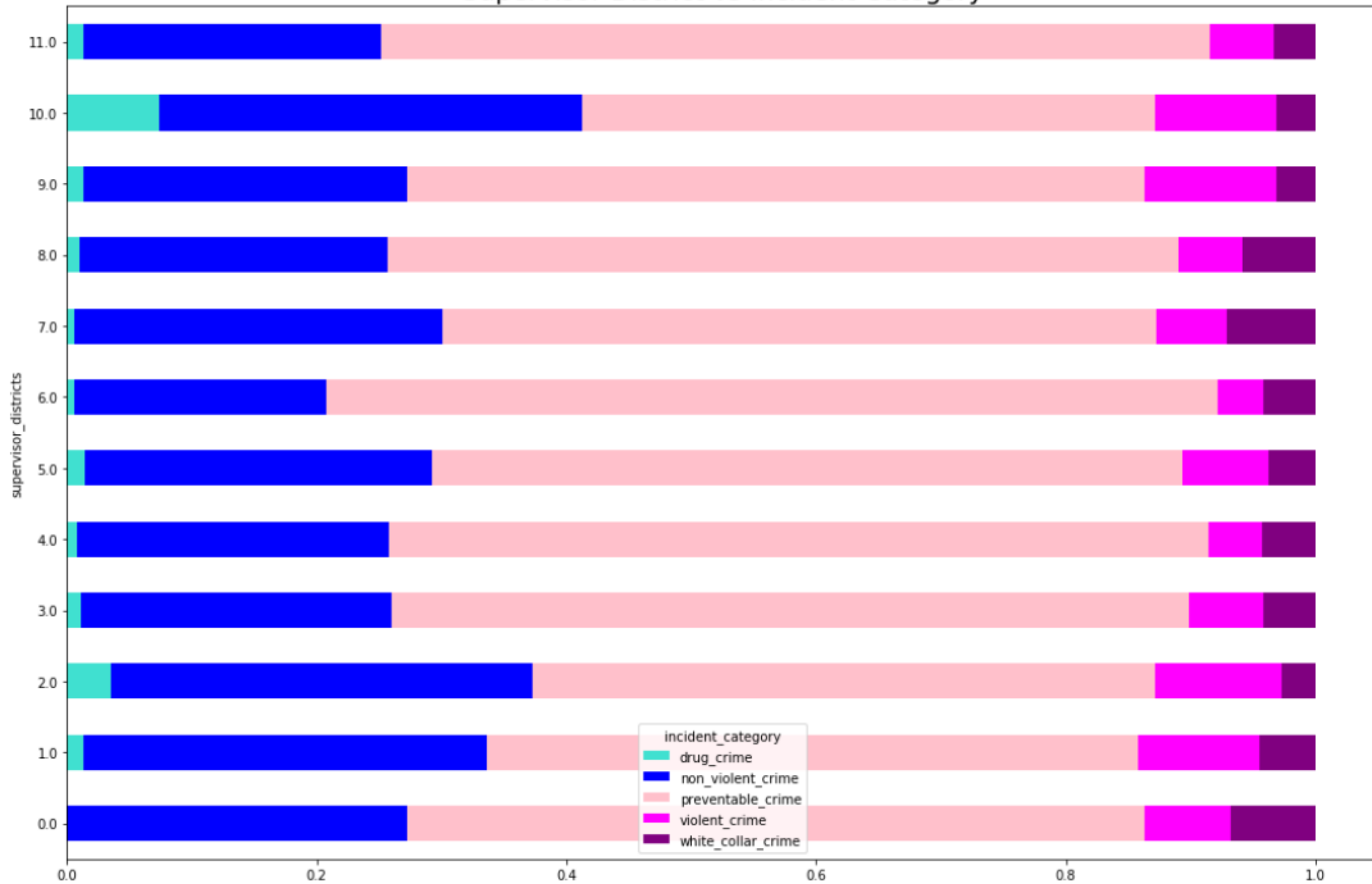
Preventable crimes were the most frequently committed category of crime



Police District vs Incident Category



Supervisor District vs Incident Category

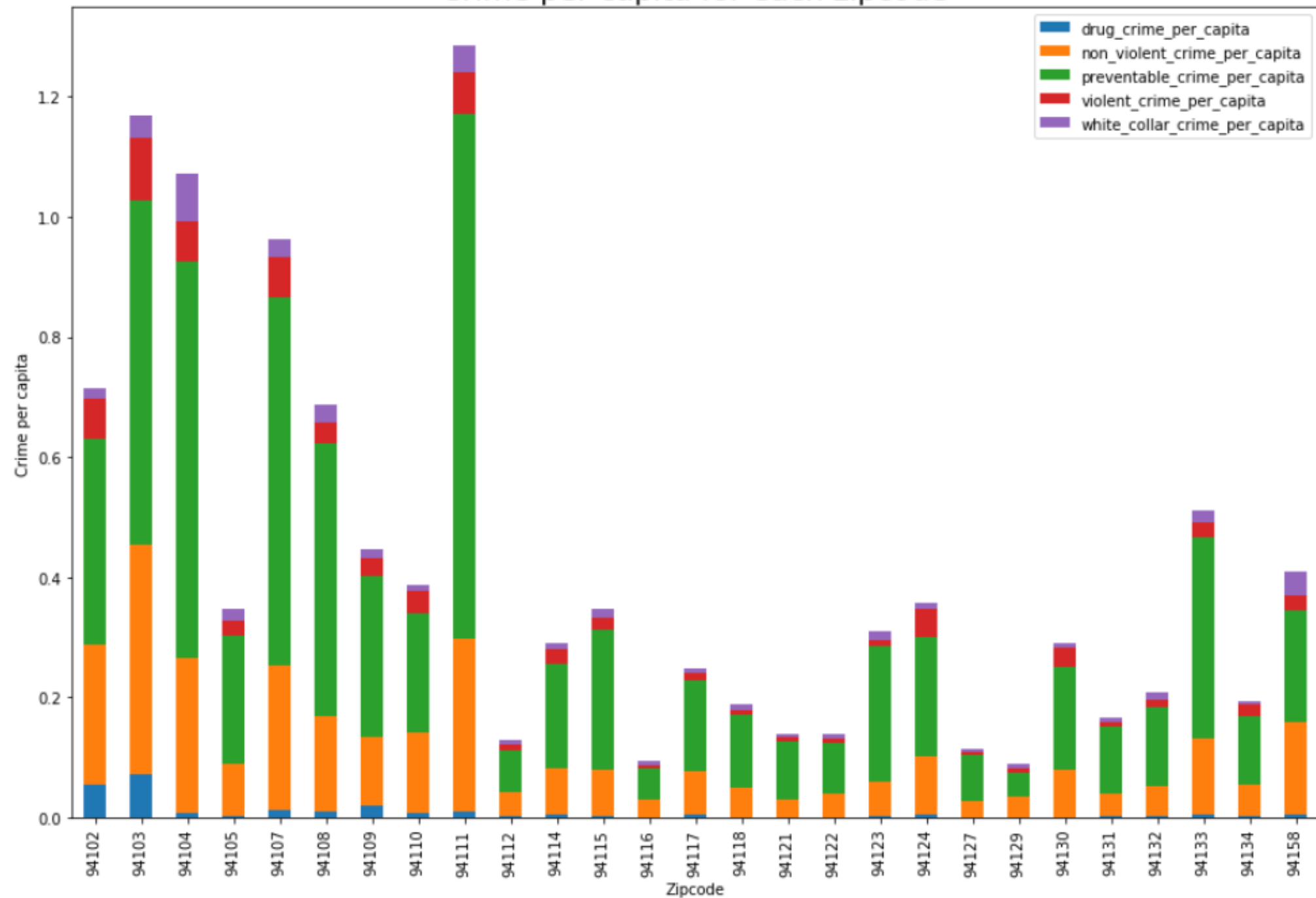


But what about the number of crimes?

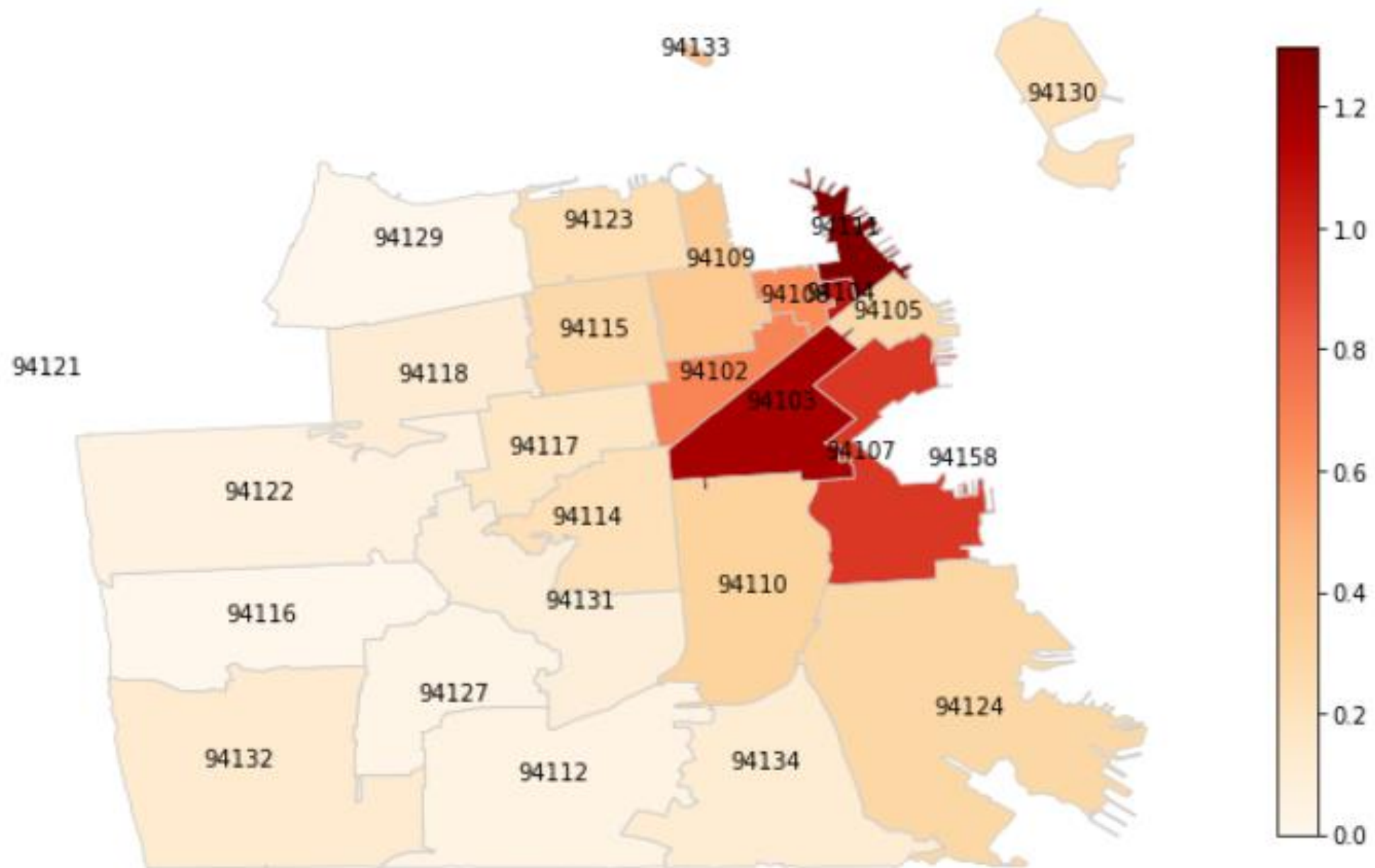
- How can we meaningfully compare the number, and not just the proportion, of crimes across different districts?



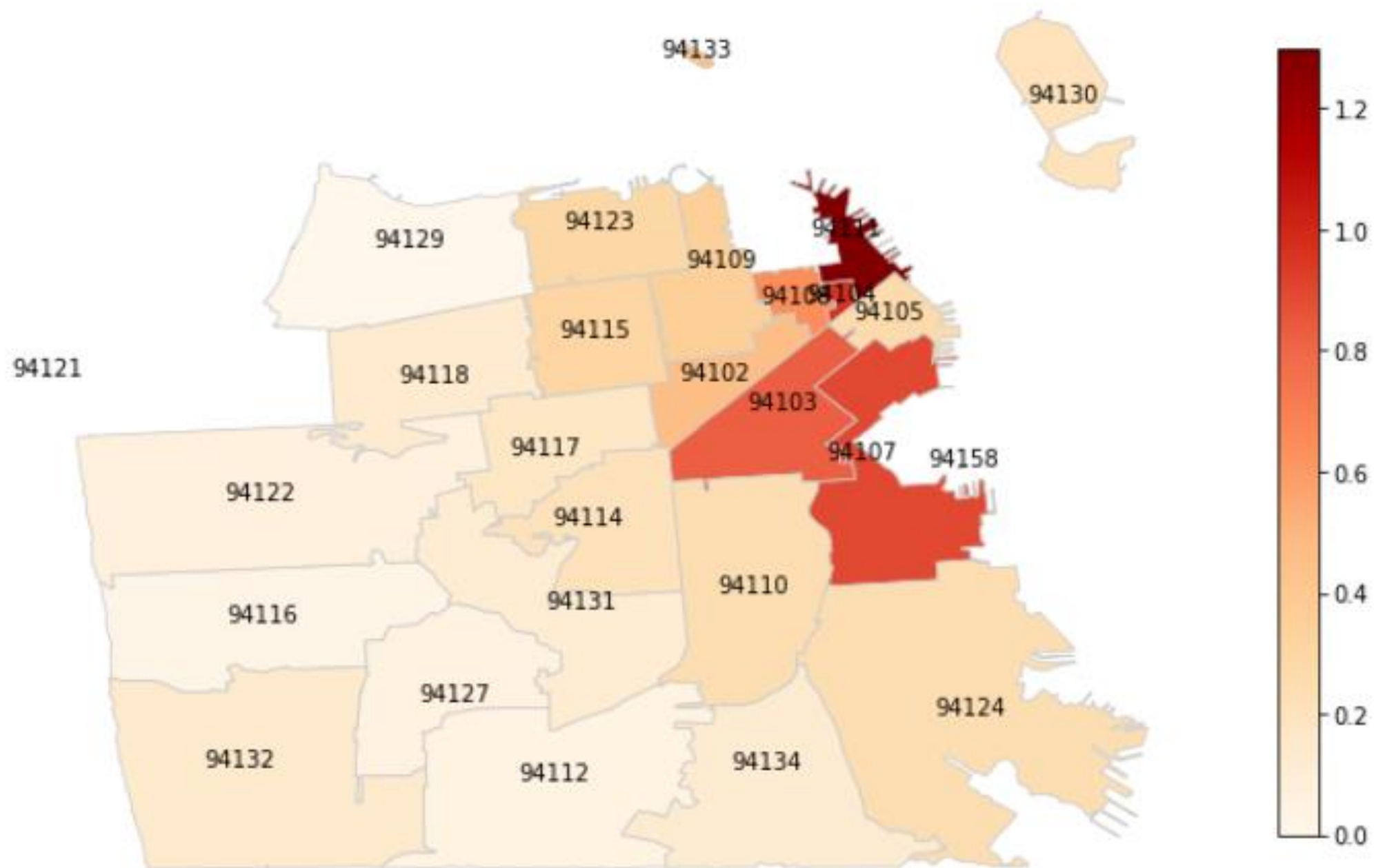
Crime per capita for each zipcode



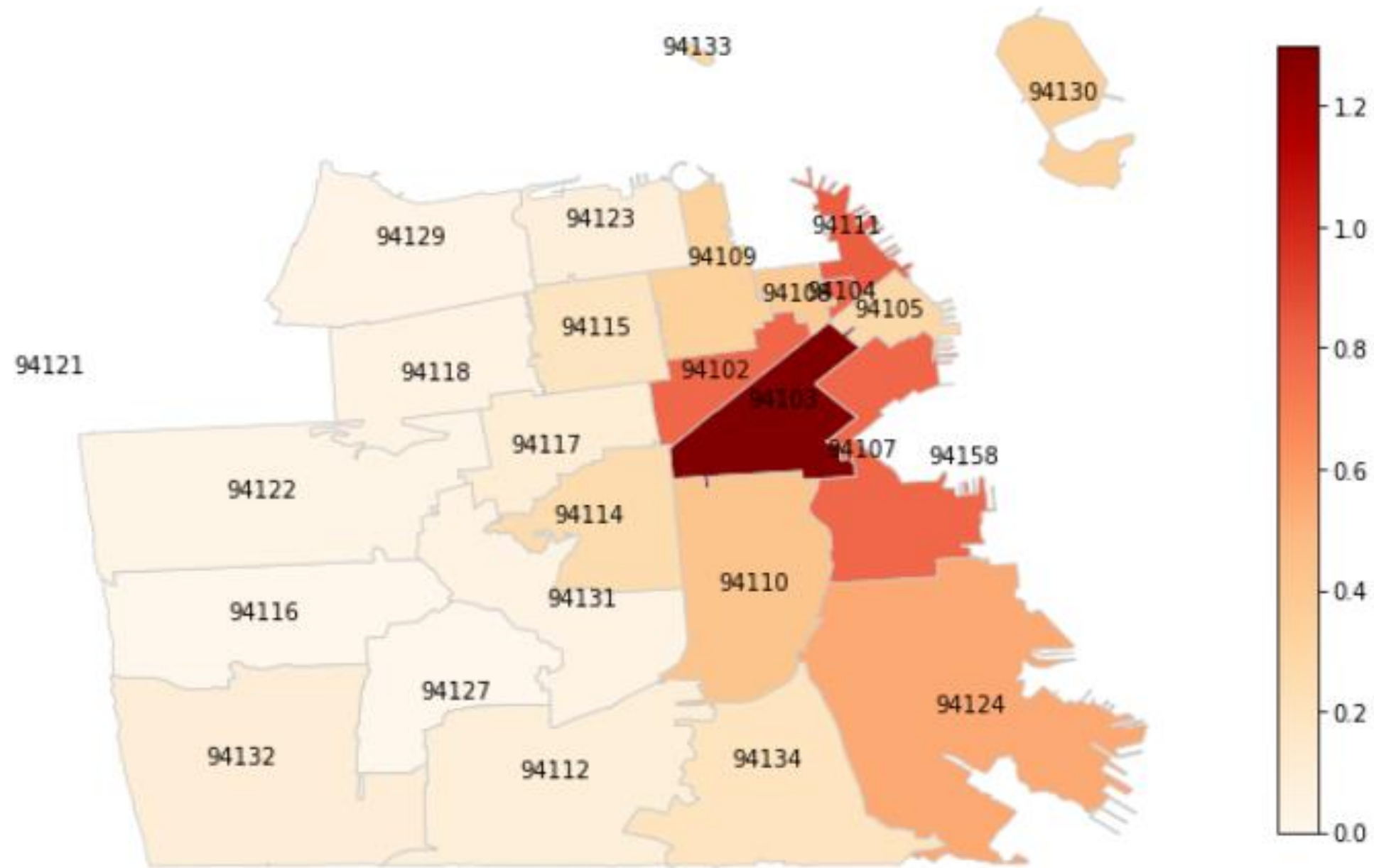
total_crime_per_capita



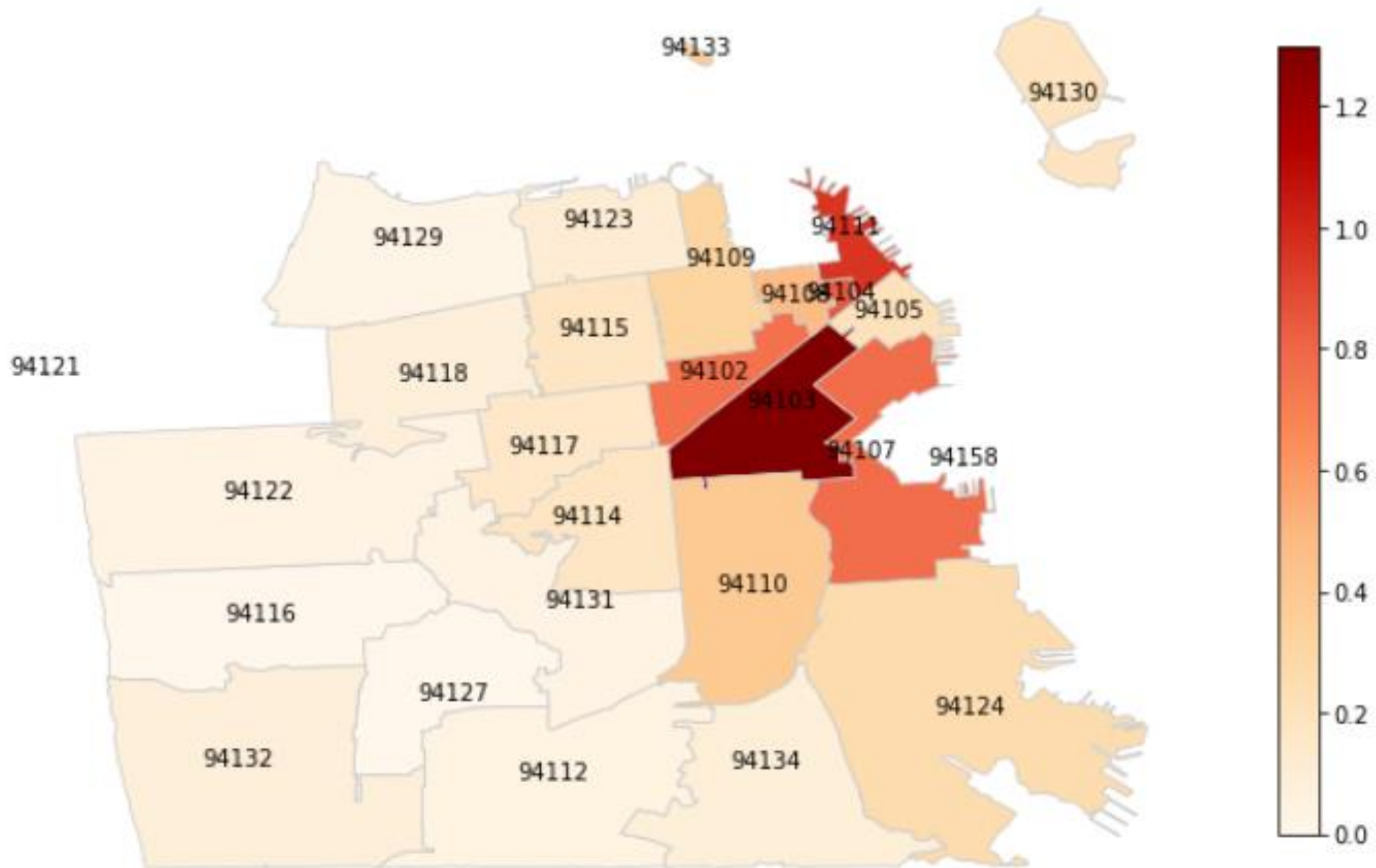
preventable_crime_per_capita



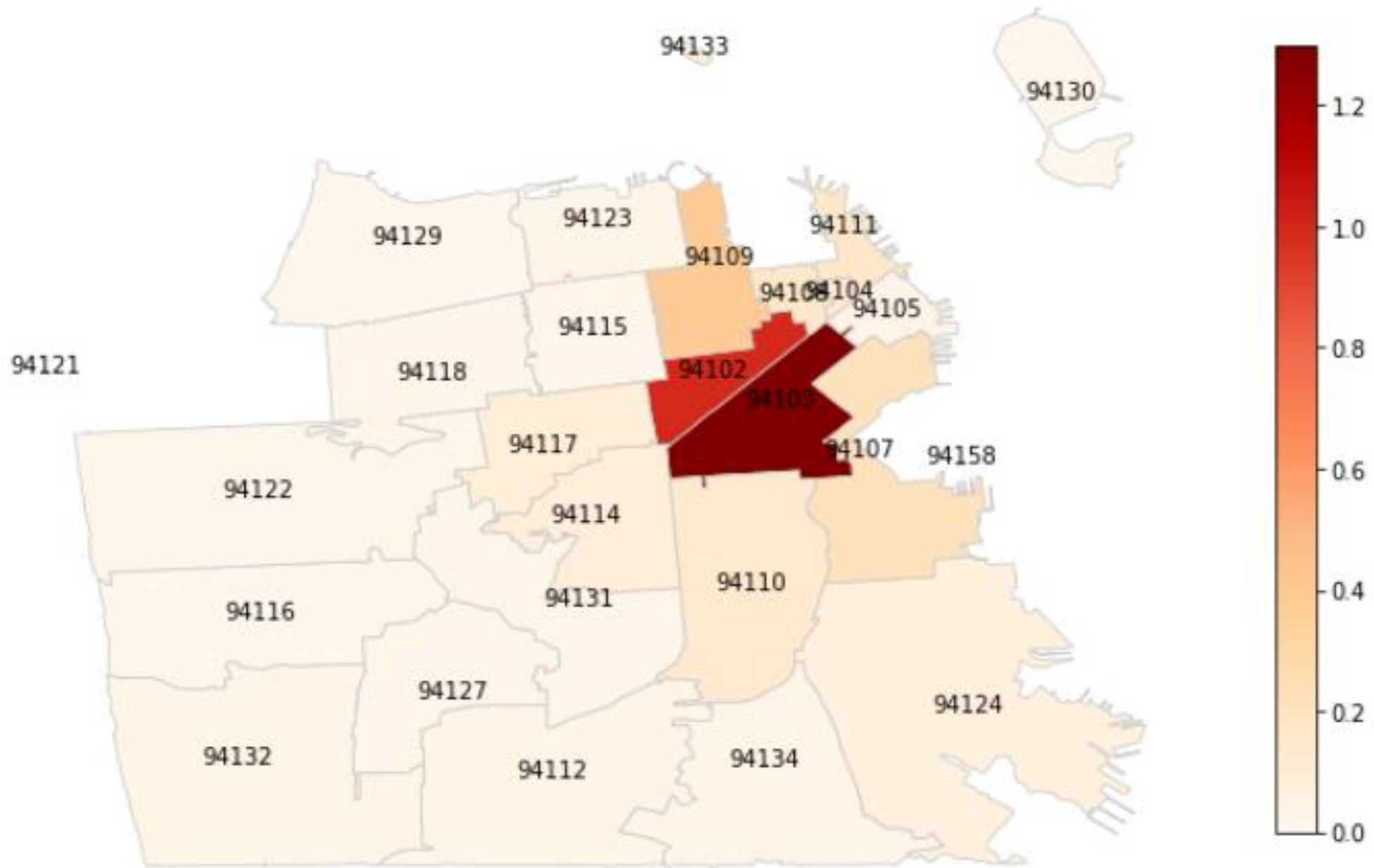
violent_crime_per_capita



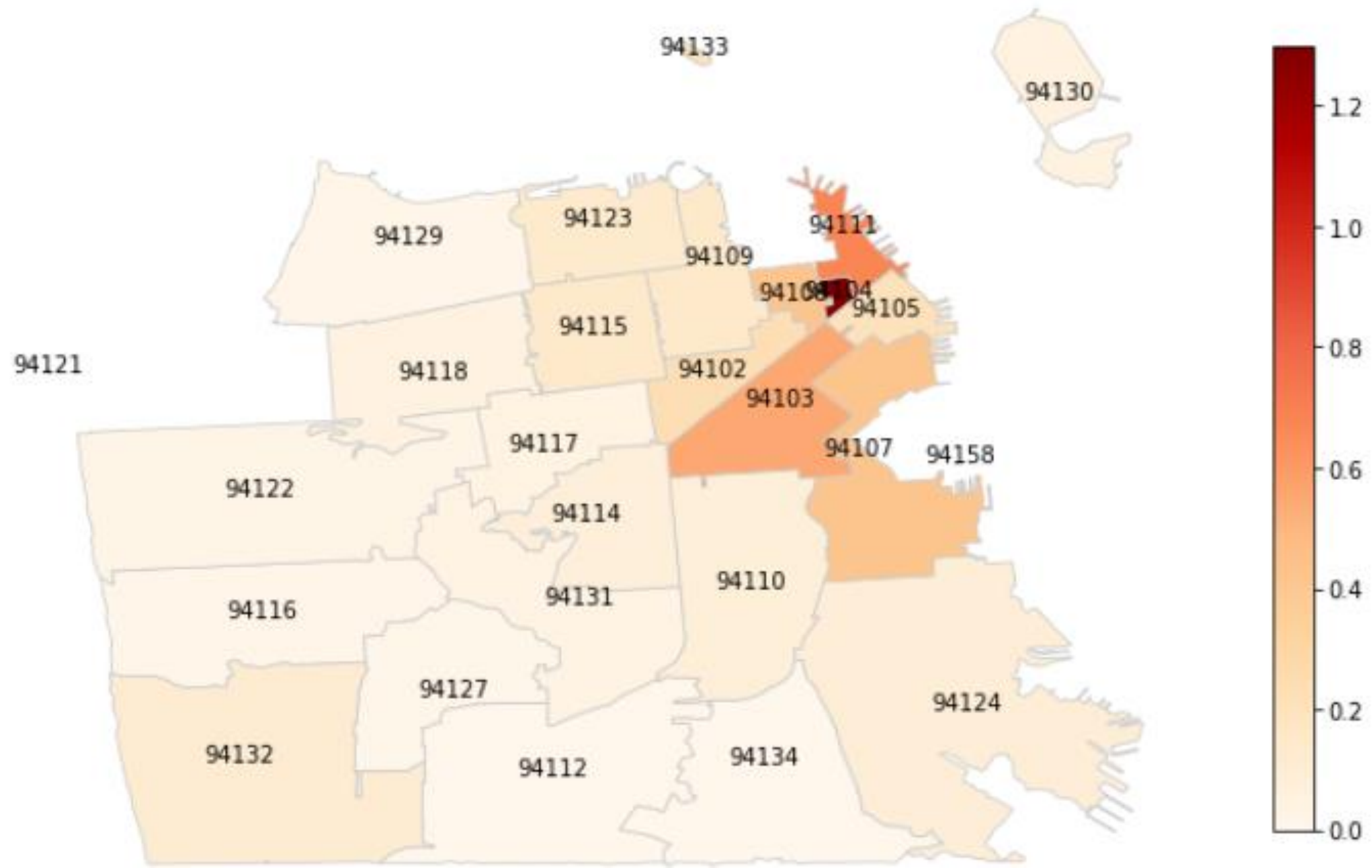
non_violent_crime_per_capita



drug_crime_per_capita



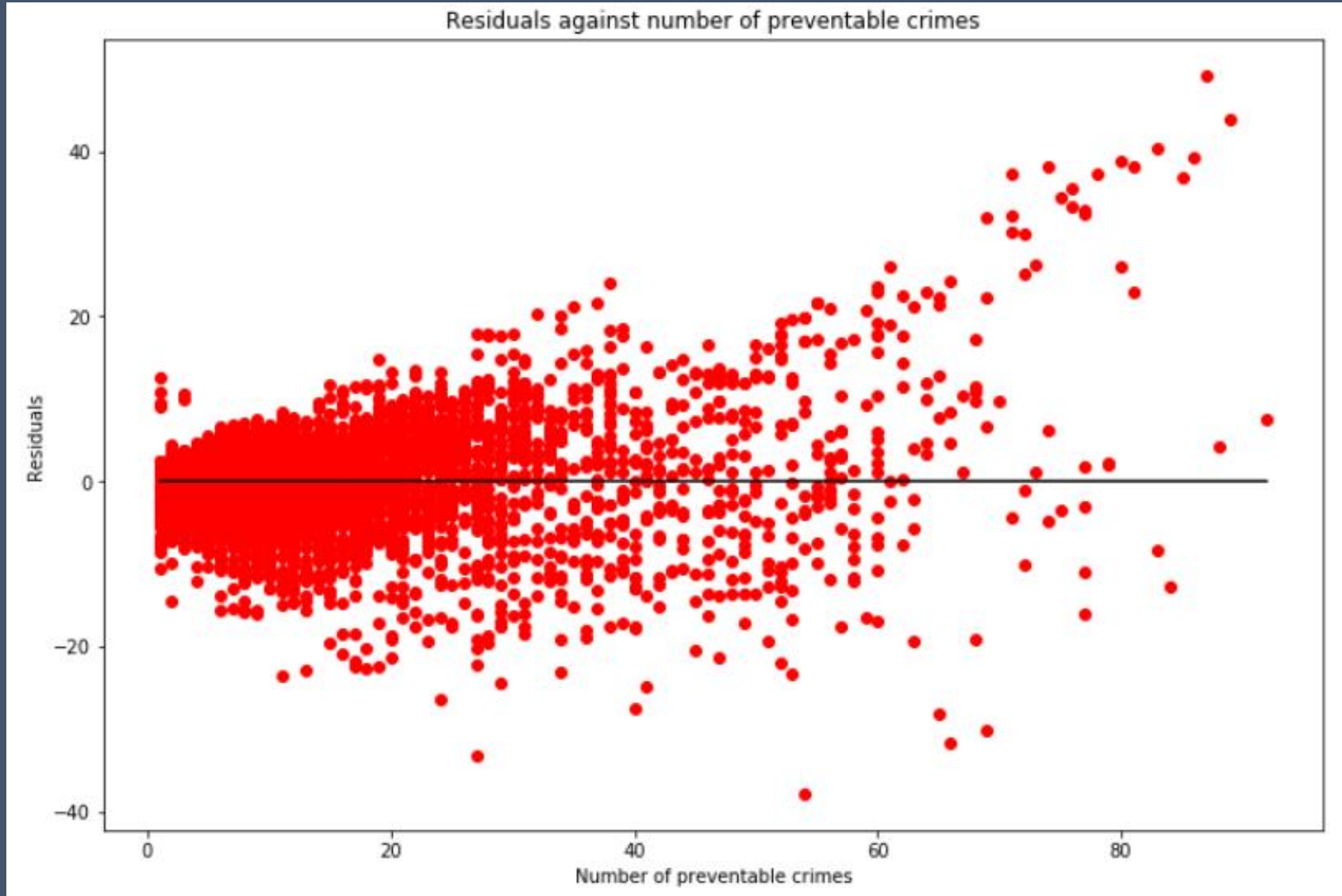
white_collar_crime_per_capita



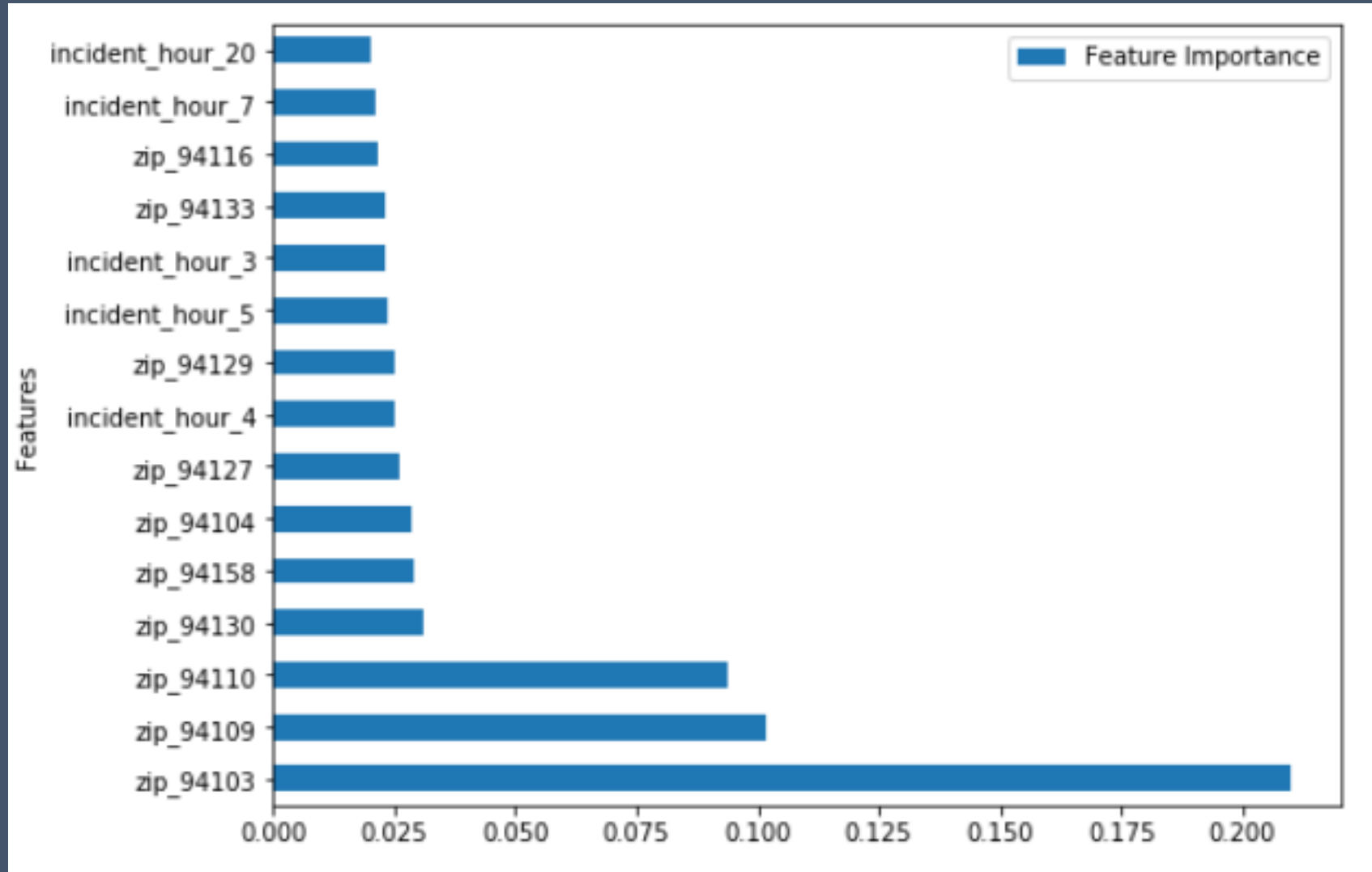
Predicting preventable crimes based on zipcode, day of week, and hour of day

Model trained on 2018 data	Mean Squared Error for X_test, y_test	R ² score for X_test, y_test
Baseline	247	-
Linear Regression	65.9	0.733
LassoCV	65.7	0.734
RidgeCV	65.8	0.734
Random Forest Regressor with GridsearchCV	47.1	0.809
XGBoost with GridsearchCV	38.3	0.845

Success with R^2 score of 0.8 when XGBoost model is used to predict from test data



Zipcodes were the most important features in predicting the number of preventable crimes



Practical Steps for SF Police to take

	incident_day_of_week	zip	incident_hour	preventable_crime	xgb_predict
2470	Thursday	94103	19	1.326923	1.907094
1260	Saturday	94103	19	1.269231	1.876870
43	Friday	94103	19	1.615385	1.858716
3079	Tuesday	94103	19	1.250000	1.791903
1259	Saturday	94103	18	1.480769	1.787161
1871	Sunday	94103	19	1.038462	1.766315
3694	Wednesday	94103	19	1.596154	1.757343
1261	Saturday	94103	20	1.480769	1.690898
648	Monday	94103	19	1.307692	1.676706
3078	Tuesday	94103	18	1.769231	1.623973

Next Steps

- Include further demographic data
- Include more years of data
- Corroborate quantitative insights with qualitative experiences of Police officers and policymakers

Any questions?