# Investigation of San Francisco Crime Data 2003-2019

• • •

By: Garrett Yamane

# Problem Statement

- **Target Client**: San Francisco Police Department

- **Goal**: Identify important crime features and use them to develop unsupervised machine learning models for suggesting patterns that contribute to crimes in San Francisco

# Data Acquisition and Wrangling

- **Sources**:
  - CABLE mainframe 2003-2018 police reports
  - Crime Data Warehouse 2018-2019 police reports (update in 2018)
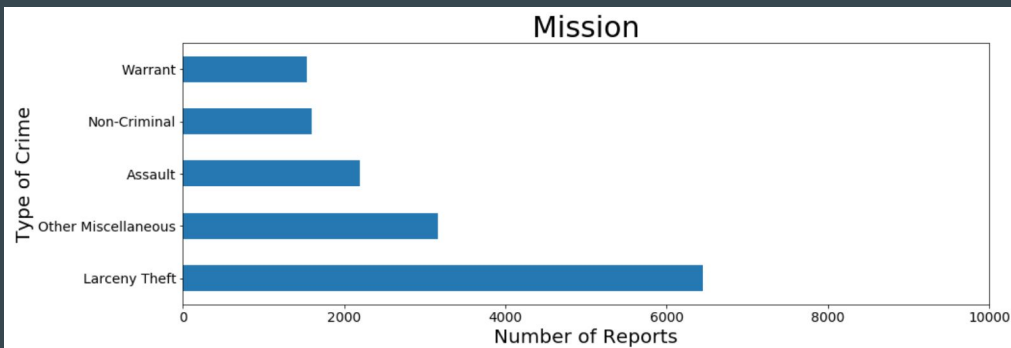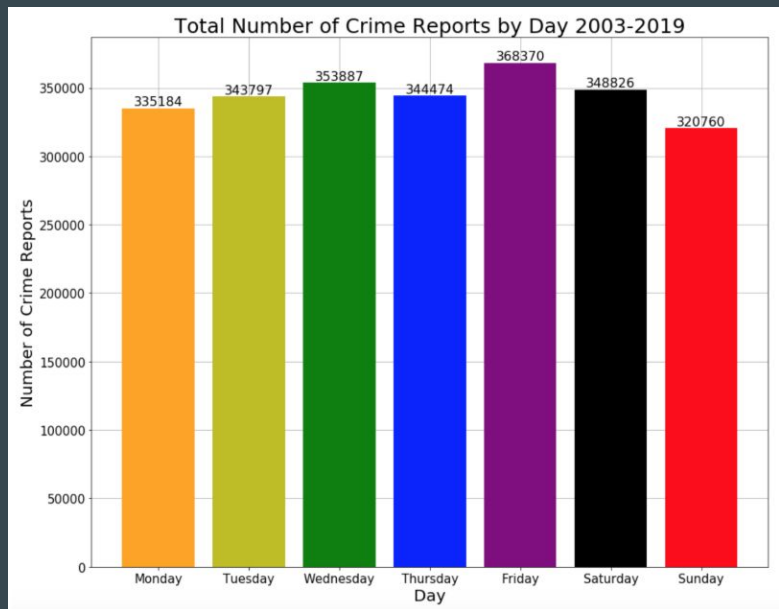
- **Final Data Frame**:
  - 2,415,298 rows (*independent crime reports*)
  - 15 crime feature columns

- **2018**: update made from CABLE mainframe to Crime Data Warehouse
        to make crime data more accessible

Crime Features for Final Data Frame
(% missing values)

| Incident Number | 0.00 |
|---|---|
| Incident Category | 0.00 |
| Incident Description | 0.00 |
| Incident Day of Week | 0.00 |
| Incident Date | 0.00 |
| Incident Time | 0.00 |
| Police District | 0.00 |
| Resolution | 0.00 |
| Intersection | 0.59 |
| Longitude | 0.59 |
| Latitude | 0.59 |
| point | 0.59 |
| Row ID | 0.00 |
| Incident Year | 0.00 |
| Analysis Neighborhood | 0.64 |

# Initial Findings



Total Number of Crime Reports by Day 2003-2019



Mission

Questions for Data Set
- Yearly crime rate change?
- What types of crimes are reported the most?
- Are crimes more likely to happen on certain days? Months? Times of the day?
- What neighborhoods are most dangerous?

# 2018-2019 Choropleth Map of Crimes per Neighborhood



| | Neighborhood | Number of Reports |
|---|---|---|
| 0 | Mission | 27168 |
| 1 | Tenderloin | 24788 |
| 2 | Financial District/South Beach | 22052 |
| 3 | South of Market | 20614 |
| 4 | Bayview Hunters Point | 13385 |
| 5 | North Beach | 7746 |
| 6 | Western Addition | 7475 |
| 7 | Castro/Upper Market | 7114 |
| 8 | Sunset/Parkside | 6935 |
| 9 | Nob Hill | 6572 |

# Application of Inferential Statistics

**Crime Report Rate**: Is there a statistical significance between the average number of crimes per year between different police districts?

- $H_0$: The true mean crime rate between the two police districts are the same

- $H_1$: The true mean crime rate between the two police districts are *not* the same

| Police District | Mean Per Year | Variance | Total Reports |
|---|---|---|---|
| southern | 26848.50 | 1.949013e+07 | 107394 |
| mission | 20342.50 | 2.423115e+06 | 81370 |
| northern | 20130.25 | 5.389582e+05 | 80521 |
| central | 20011.50 | 1.041427e+07 | 80046 |
| bayview | 13818.50 | 1.223255e+06 | 55274 |
| ingleside | 11874.50 | 1.143297e+06 | 47498 |
| taraval | 11501.00 | 5.272353e+05 | 46004 |
| tenderloin | 11223.75 | 8.588913e+06 | 44895 |
| richmond | 8993.25 | 1.796556e+05 | 35973 |
| park | 8479.75 | 7.493216e+05 | 33919 |

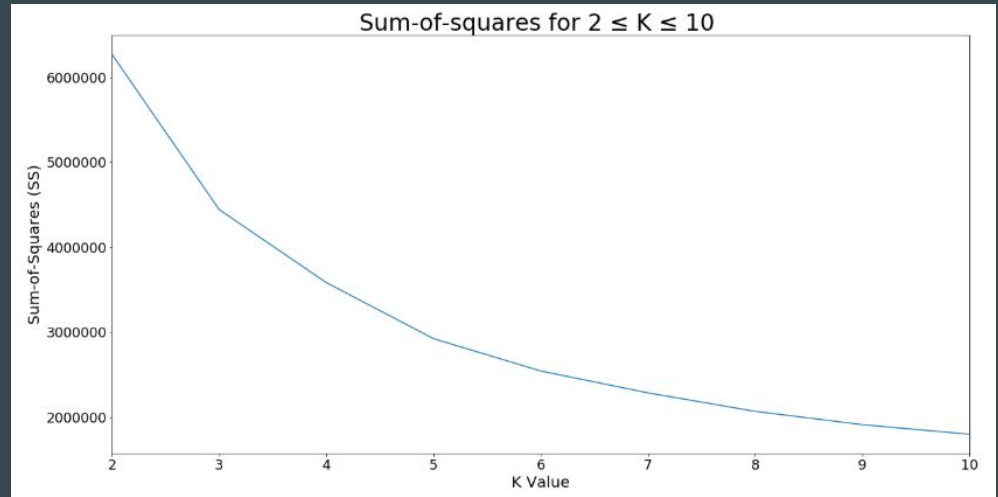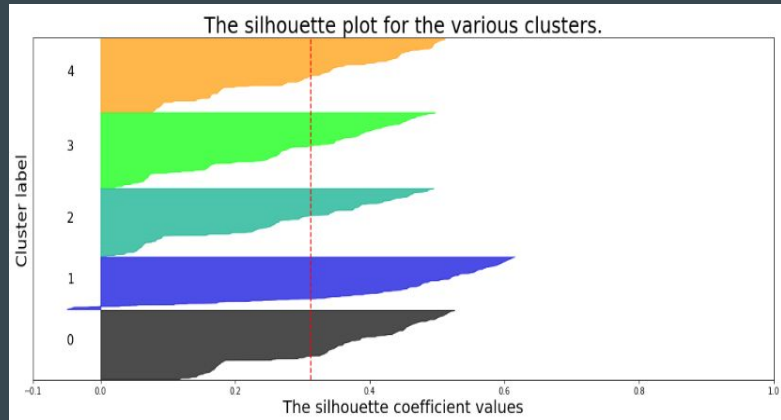# Application of Inferential Statistics: Feature Selection

- **Chi-Square Test**
  - **Target Feature:** *Incident Category*
  - **Test Features:**
    - *Incident Hour*
    - *Incident Day of Week*
    - *Incident Month*
    - *Police District*
    - *Resolution*
- **Goal**
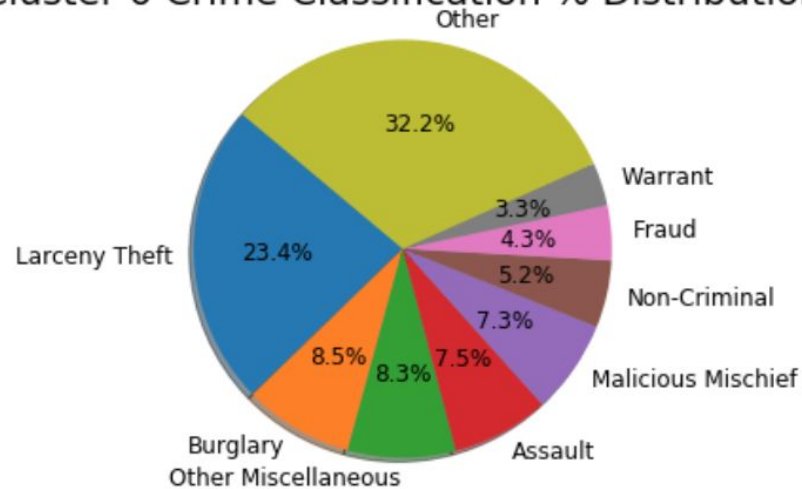  - What features are useful determining the type of crime to be committed?

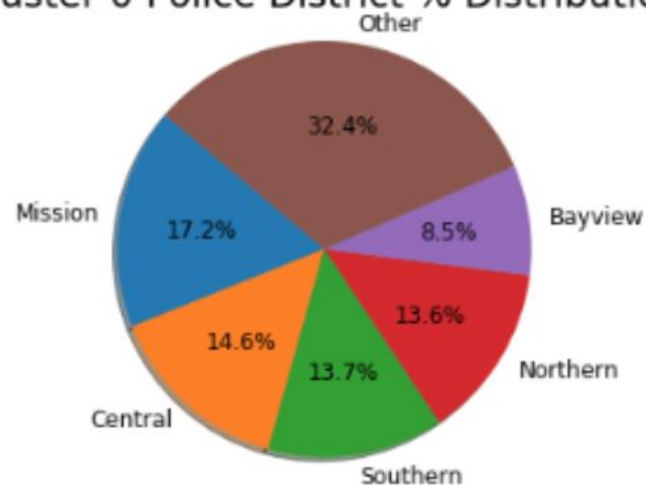# Baseline Clustering Model: K-Means



The silhouette plot for the various clusters.



Sum-of-squares for 2 ≤ K ≤ 10

*Baseline k-means model built with k = 5
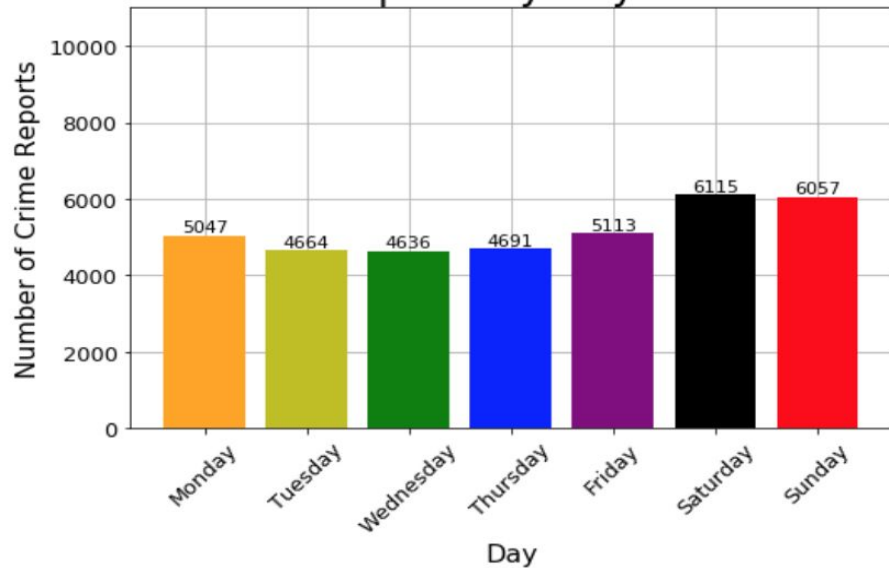*5 features used to cluster: Incident Hour, Incident Day of Week, Incident Month, Police District, Incident Category
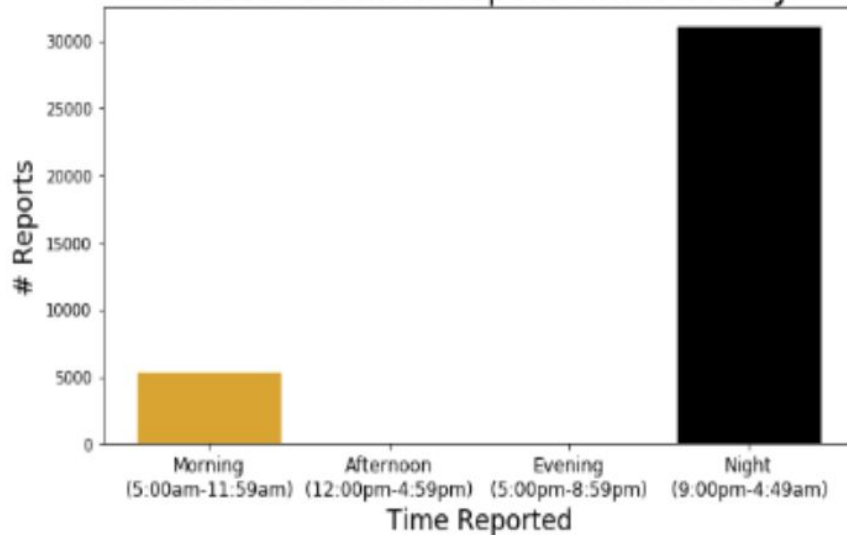
Cluster 0 Crime Classification % Distribution

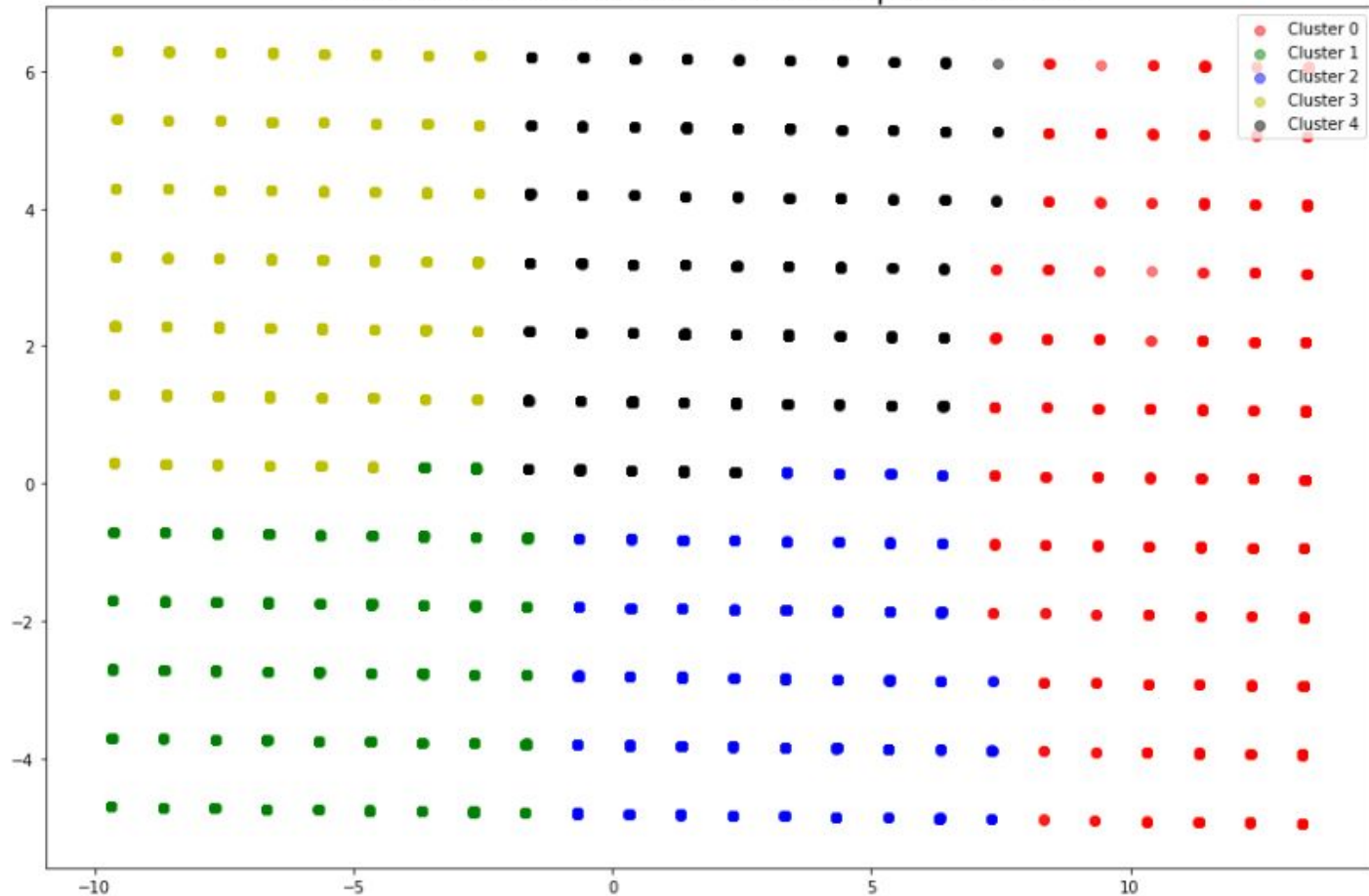Cluster 0 Police District % Distribution

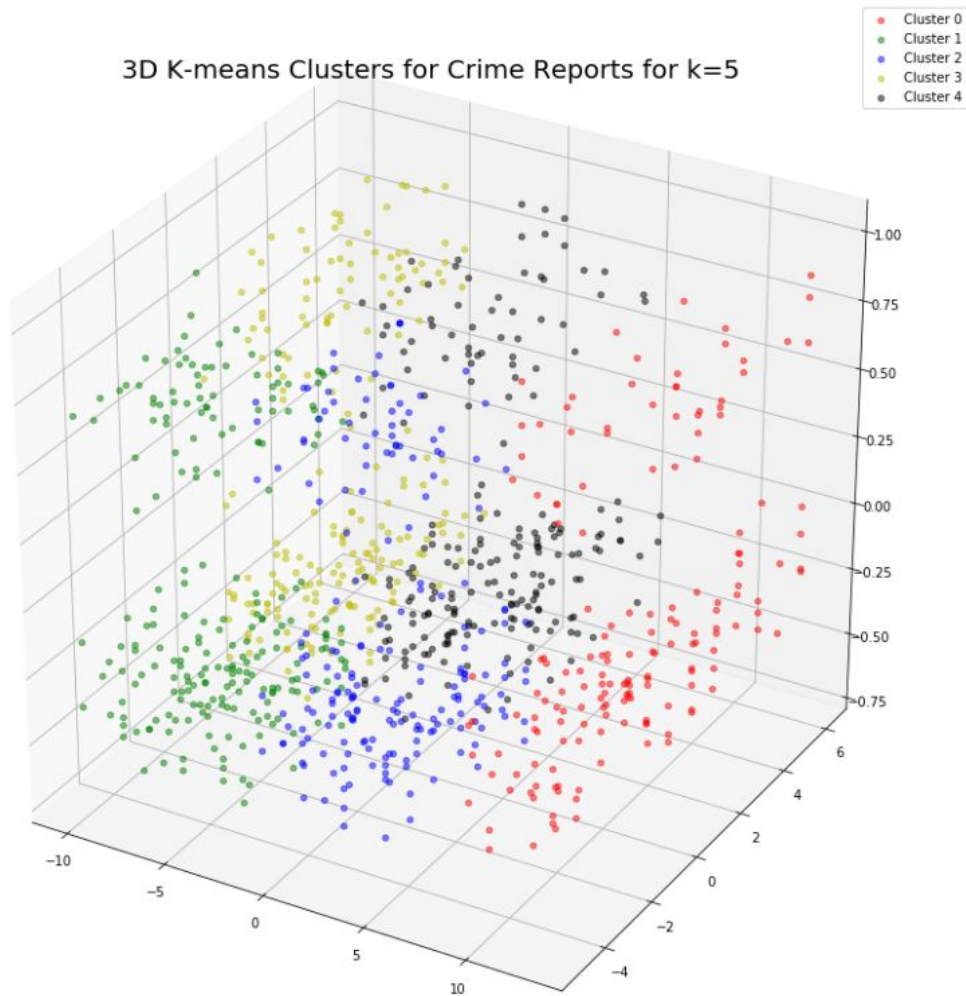Crime Reports by Day Cluster 0

| Day | Number of Crime Reports |
|---|---|
| Monday | 5047 |
| Tuesday | 4664 |
| Wednesday | 4636 |
| Thursday | 4691 |
| Friday | 5113 |
| Saturday | 6115 |
| Sunday | 6057 |

Cluster 0 Crime Report Times of Day

| Time Reported | # Reports |
|---|---|
| Morning (5:00am-11:59am) | ~5000 |
| Afternoon (12:00pm-4:59pm) | ~0 |
| Evening (5:00pm-8:59pm) | ~0 |
| Night (9:00pm-4:49am) | ~31000 |

# Visualizing K-Means Clusters: PCA

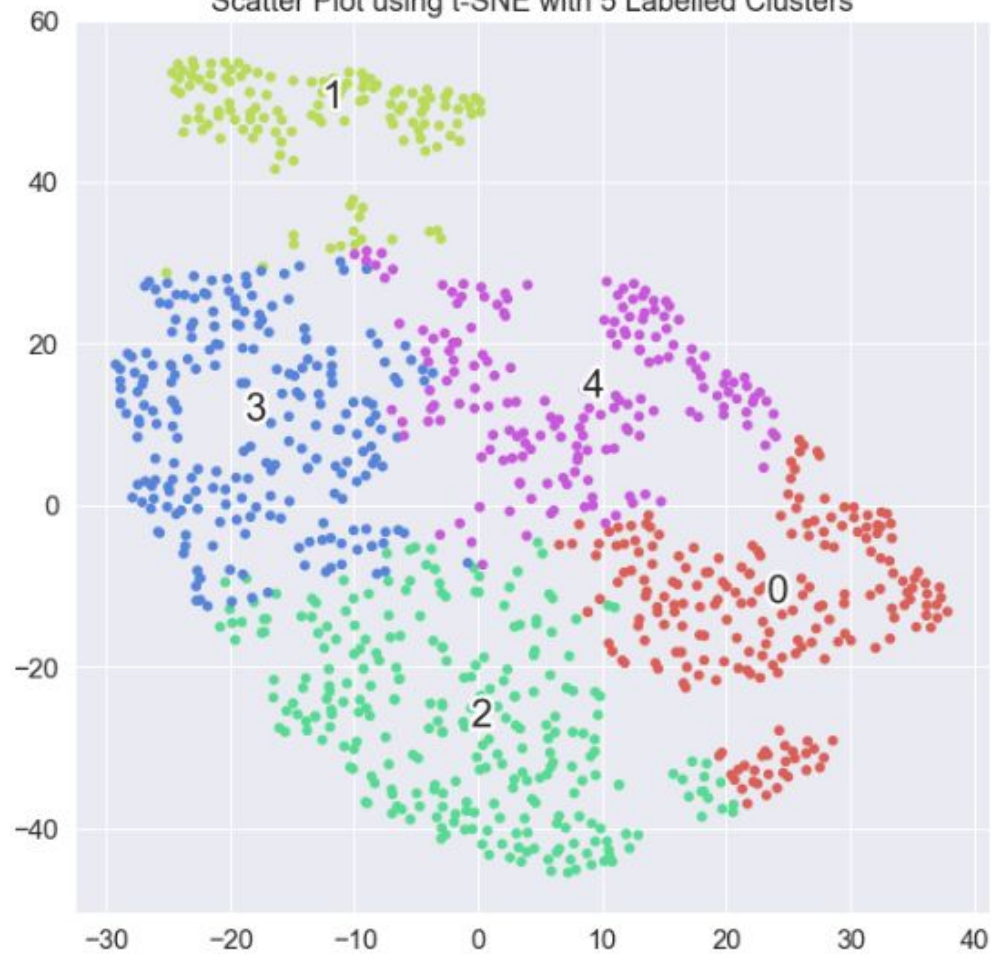2D K-means Clusters for Crime Reports for k=5

3D K-means Clusters for Crime Reports for k=5

# Model Extension: Application of t-distributed stochastic neighbor embedding (t-SNE)

Scatter Plot using t-SNE with 5 Labelled Clusters
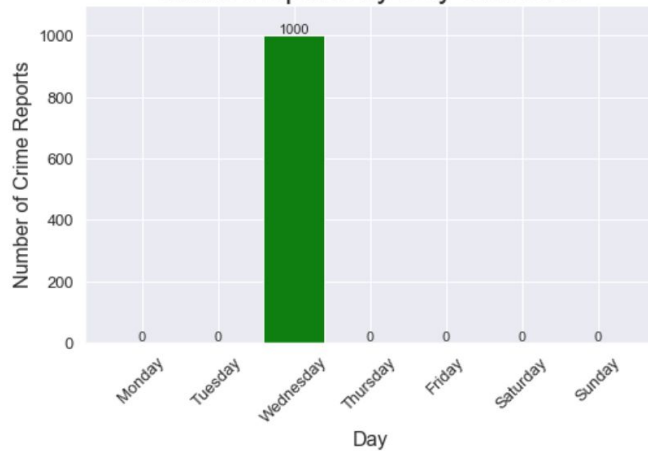
# Model Extension: K-Modes

- Better for handling categorical data
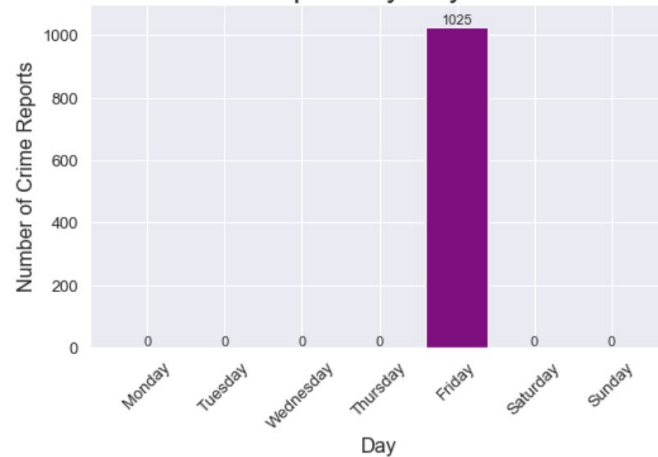- Distance Function: dissimilarity between objects rather than Euclidean distance



Cost for K-Modes Models 2 ≤ K ≤ 10

*Used same features to cluster as k-means model*

# K-Modes: Key Differences

# Conclusions and Future Work

- **Big takeaways**
  - No single algorithm can reveal all the underlying patterns of my crime data
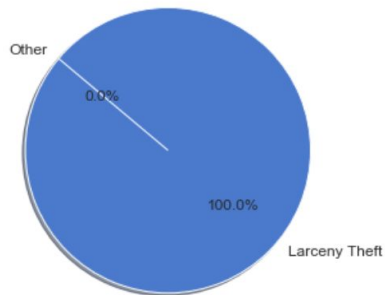  - K-modes specifically targets categorical data -> stronger representation of SF crime data

- **Future Work**
  - Configure own crime classifications rather than dealing with every unique category
  - Supervised Learning Perspective
    - Predict the type of crime that will happen given certain conditions (time of day, location, etc.)

# Recommendations to Client

- "Larceny Theft" by far the most reported crime
  - Large impact on crimes that could be connected to each other
- Target police districts such as Central and Northern with high priority
- Allocate more resources to districts where crimes are most likely may help reduce and prevent future crimes