

# Raj Ramrup

Data: NYPD Complaint Data 2006- October 2017

Goal: To find the risk of being at a certain place and time.

Machine learning problem:

If a crime is to occur at a certain time and place, what type of crime will it be (assault, robbery, etc.)

Given: Date, Time, Location

Predict: Crime Type ( Classification)

Limitations: False complaints, crimes that go unreported

# Cleaning Features

- Take the features from the dataset
- Divide them into parts
  - Date → weekday, year, month, day, is\_month\_start, is\_month\_end, is\_quarter\_start, is\_quarter\_end, is\_year\_start, is\_year\_end
  - Time → hour, minute, second, *is\_business\_hours*
  - Latitude, Longitude →  
 $x = \text{np.cos(df['Latitude'])} * \text{np.cos(df['Longitude'])}$  )  
 $y = \text{np.cos(df['Latitude'])} * \text{np.sin(df['Longitude'])}$  )  
 $z = \text{np.sin(df['Latitude'])}$  )
  - Typos:
    - Things like '1017' as a year

# Data

- Data for 11 years: 5 million rows
- Data for this year: 300,000 rows
- SVM classification
  - ~20 Features
  - ~80 Classes
  - 50/50 split between training and test

# Results?

- Inconclusive...

KY_CD	Three digit offense classification code
OFNS_DESC	Description of offense corresponding with key code
PD_CD	Three digit internal classification code (more granular than Key Code)

- Currently run on 20% of only this years data. Bad accuracy.
- Attempted to run on whole dataset (with bad feature) but that took over 14 hours and never finished
- If run on the whole dataset:
  - Monthly trends
  - Seasonal trends
  - More data on area trends

# Uses

- Take the locations of subway terminals and various times
- → see which train lines are more dangerous
- → find the least dangerous route home at 3 am