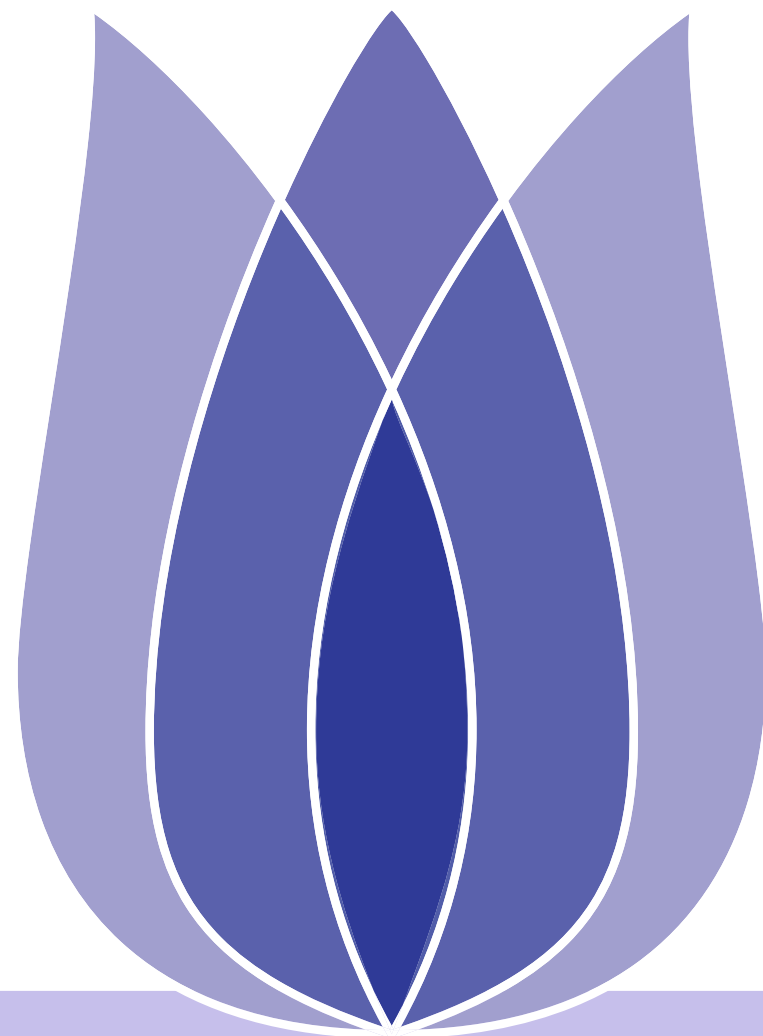# San Francisco Crime classification

Jia Huang

Xi'an Shiyou University

Chinese Academy of Sciences

October 16, 2020

# Project Overview

# Project Background And Purpose

■ Background

From 1934 to 1963, San Francisco was infamous for housing some of the world's most notorious criminals on the inescapable island of Alcatraz. Today, the city is known more for its tech scene than its criminal past. But, with rising wealth inequality, housing shortages, and a proliferation of expensive digital toys riding BART to work, there is no scarcity of crime in the city by the bay. From Sunset to SOMA, and Marina to Excelsior, this dataset provides nearly 12 years of crime reports from across all of San Francisco's neighborhoods.

■ Purpose

predict the category of crime that occurred, given the time and location visualize the city and crimes (see Mapping and Visualizing Violent Crime for inspiration) Content.

Defn

# Data Pre-Processing

# Feature Item

This dataset contains incidents derived from SFPD Crime Incident Reporting system. The data ranges from 1/1/2003 to 5/13/2015. The training set and test set rotate every week, meaning week 1,3,5,7... belong to test set, week 2,4,6,8 belong to training set. There are 9 variables:

- Characteristic Term
  - Dates
  - Category
  - Descript
  - DayOfWeek

- Characteristic Term
  - PdDistrict
  - Resolution
  - Address
  - X
  - Y

## Feature Item

By making a comprehensive analysis of the nine characteristic items in the data set mentioned above, we can reach the following conclusions:

```
First date:  2003-01-06 00:01:00
Last date:  2015-05-13 23:53:00
Test data shape   (878049, 9)
```

Figure 1

The data range was from 1/1/2003 to 5/13/2015, and a data training set containing nine feature items and 87,8049 samples was created.

# Feature Analysis

- Statistics Were Made By Type Of 'Year' And 'Month'
  Based on a comprehensive analysis of the data set provided by Kaggle's website, it is clear that there are fewer crimes in summer and winter than in spring and fall.Therefore, a "seasonal" feature column can be added to the feature analysis.

- By 'DayOfWeek' And 'Hour' Type
  Friday saw the highest number of crimes, probably because of the American tradition of Friday parties.Sunday has the lowest crime rate.So you can add the "weekend or not" feature column.Crime was lowest in the early hours of the morning and highest at 12 o 'clock and 17 and 18 o 'clock in the evening.Therefore, the time zone can be divided and the "time zone" feature column can be added

# Feature Selection

# Modelling

# Calculate the Baseline Value For The Model

Since this is a typical multi-classification problem, we can choose to use many kinds of algorithms, including naive Bayes, KNN, decision tree and random forest.

# Model Optimization

# Ideas Improvement