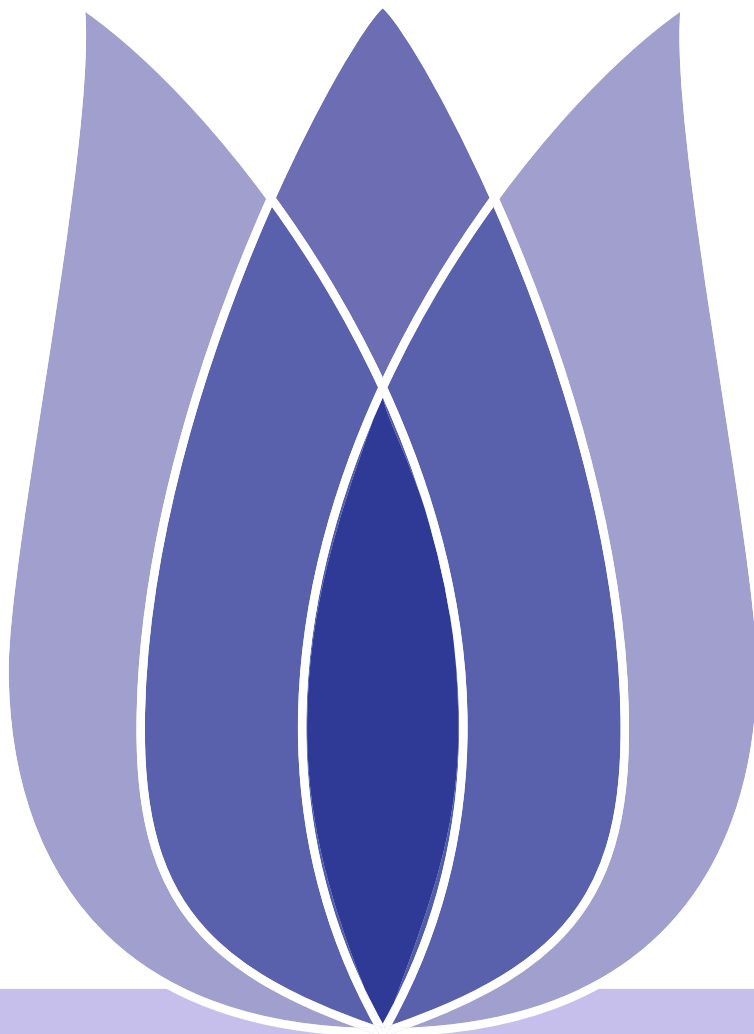


San Francisco Crime classification

Jia Huang

Xi'an Shiyou University
Chinese Academy of Sciences

September 27, 2020





Overview



Project Overview



TULIP

Team for Universal Learning and Intelligent Processing



Project Background And Purpose

Defn	<div><div>■ Background</div><div>From 1934 to 1963, San Francisco was infamous for housing some of the world’s most notorious criminals on the inescapable island of Alcatraz. To-day, 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.</div><div>■ Purpose</div><div>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.</div></div>
------	---



Data Pre-Processing



TULIP

Team for Universal Learning and Intelligent 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



Dates - timestamp of the crime incident

Category - category of the crime incident (only in train.csv). This is the target variable you are going to predict.

Descript - detailed description of the crime incident (only in train.csv)

DayOfWeek - the day of the week

PdDistrict - name of the Police Department District

Resolution - how the crime incident was resolved (only in train.csv)

Address - the approximate street address of the crime incident

X - Longitude

Y - Latitude





Feature Analysis



TULIP

Team for Universal Learning and Intelligent Processing



- 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



TULIP

Team for Universal Learning and Intelligent Processing



Modelling



TULIP

Team for Universal Learning and Intelligent Processing



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



TULIP

Team for Universal Learning and Intelligent Processing



Ideas Improvement



TULIP

Team for Universal Learning and Intelligent Processing