# CPSC 483
# Data Mining and Pattern Recognition



**Department of Computer Science**

**California State University, Fullerton**
**Spring 2016**

# Group Project

**Submitted by:**

| Name | Email ID | CWID |
|------|----------|------|
| Anurag Patsariya | anurag.patsariya@csu.fullerton.edu | 803000082 |
| Nitesh Joshi | niteshjoshi@csu.fullerton.edu | 802385211 |
| Vidhi Patel | vidhi.patel38@csu.fullerton.edu | 893375105 |

# 1. Introduction

Kaggle competition is a data science competition where different companies put their data for predication and analysis and data miners from all over the world participate and produce best models using the given data. In Kaggle, we choose a particular competition and we download the training data for the same. Next step is, building a model using our choice of methods or tools. Finally, the predications are uploaded and Kaggle scores the solution, which are then displayed on the leaderboard.
The Kaggle competitions can be found on this url: https://www.kaggle.com/competitions

We had to choose a competition which was active on Kaggle. Therefore, we looked at different competitions and found one which was on San Francisco Crime Classification: https://www.kaggle.com/c/sf-crime

This dataset is brought by SF Open Data, the central clearinghouse for data published by the City and County of San Francisco. 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.

## 1.1 Environment:
Python -3.5

## 1.2 Libraries:
1. **Pandas: -** It offers data structures and operations for manipulating numerical tables and time series
2. **Sklearn: -** They have all algorithm with their functions.
3. **Numpy: -** It helps in mathematical computation.
4. **Matplotlib: -** This library is used to plot graphs in python, it has inbuilt functions for it.
5. **Math:** - This library is used to access the mathematical function.

## 1.3 Tools Used:
1. Visual Studio Code.
2. RapidMiner

# 2. Installation of Python and it libraries

**Step 1:**
Download python: - **https://www.python.org/downloads/**
**Step 2:**
Install all the libraries needed:
1. Install pip(This will help us install other libraries)
2. Install numpy (pip install numpy)
3. Install sklearn (pip install sklearn)
4. Install pandas (pip install pandas)
5. Install matplotlib (pip install matplotlib)
6. Install math (pip install math)

# 3. Download data from Kaggle

**Data:** https://www.kaggle.com/c/sf-crime/data
- Download all the files on this link.
- Train.csv: - this data will help us to model the data using the algorithm.
- Test.csv: - the modeled algorithm will help us to predict on the test data.
- samplesubmission.csv: - the predicted values need to be stored in the format shown in this file.

**Data: -**

- Dates - timestamp of thse 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

# 4. Preprocessing of the Data
- **Date:** Data we got from kaggle has date in String format. So we used parse_dates to convert into it in date time format. Then we used get_dummies from Pandas library

to extract hour.

- **Categories:** We used label_encoder to convert crime categories to unique integer values.
- **Day of week:** We used get_dummies function from Pandas to convert   days into binarized array.
- **PdDistrict:** We used get_dummies function from Pandas to convert   districts into binarized array.

```
import pandas as pd

train = pd.read_csv("/Users/nikunjpatel/Desktop/DataMiningGroup/train.csv",parse_dates = ['Dates'])
test = pd.read_csv("/Users/nikunjpatel/Desktop/DataMiningGroup/test.csv", parse_dates = ['Dates'])

le_crime = preprocessing.LabelEncoder()
crime = le_crime.fit_transform(train.Category)

days = pd.get_dummies(test.DayOfWeek)
district = pd.get_dummies(test.PdDistrict)

hour = test.Dates.dt.hour
hour = pd.get_dummies(hour)
test_data = pd.concat([hour, days, district], axis=1)
```

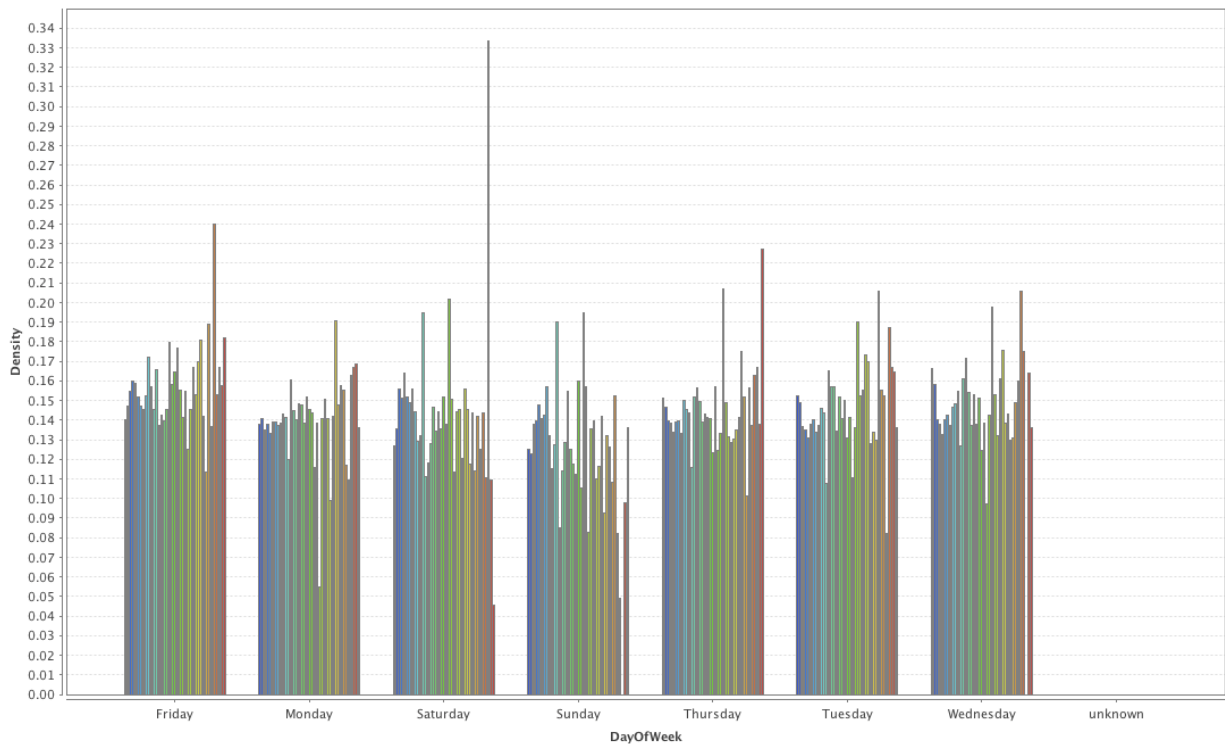# 5. Analysis of the data (using RapidMiner)

Here are some graphs produced using RapidMiner for analyzing the data. After analyzing the data, we can say that crime rates are high in Bayview, Southern and Mission districts.
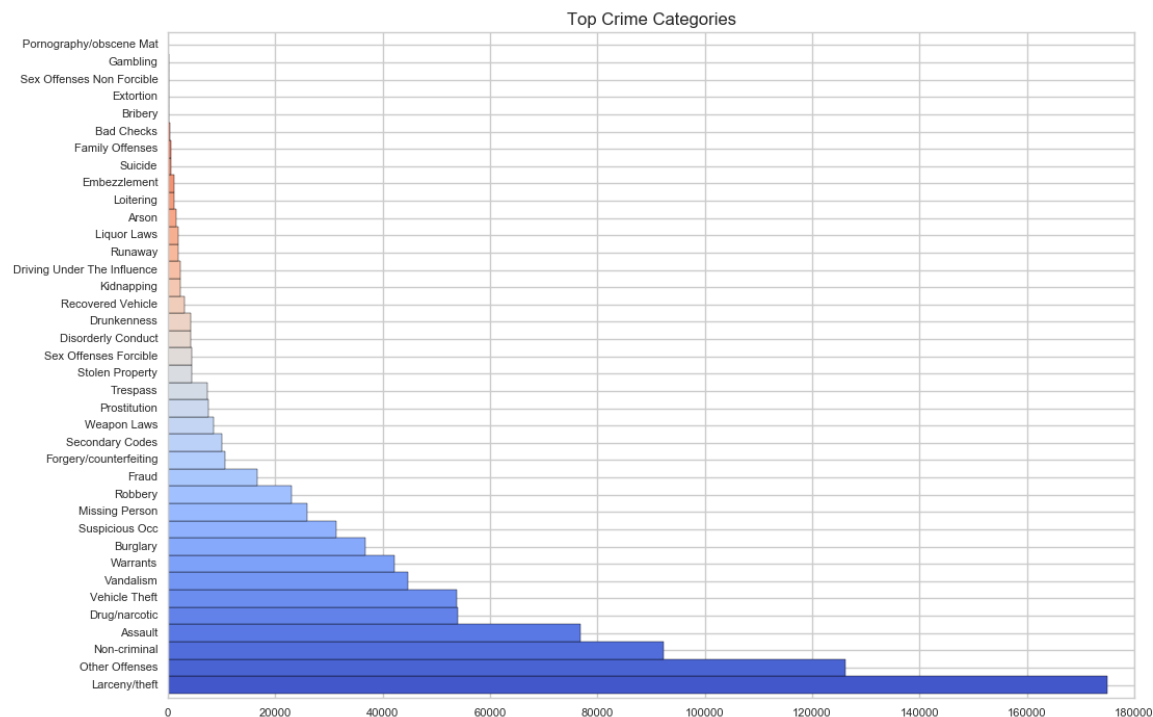
By analyzing data for the week days we can say that crime rate is high on Fridays.

Top Crime Categories

# 6. Implemented Algorithms:

### 1. K nearest neighbor algorithm:
For this algorithm we have considered
- Category
- Dates
- x and y axis.

We took these features so that we might get the exact location, but when we tried it on real algorithm it gave us some astonishing results as follows.

Output:

```
zipfile.ZipFile('/Users/nikunjpatel/Desktop/DataMiningGroup/newData.zip')
test = pd.read_csv(open('test.csv'), parse_dates=['Dates'])
x_test = test[['DayOfWeek', 'PdDistrict']]
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(x, y)
outcomes = knn.predict(x_test)
submit = pd.DataFrame({'Id': test.Id.tolist()})
for category in y.cat.categories:
    submit[category] = np.where(outcomes == category, 1, 0)
```

```python
    print("Printing")
    submit.to_csv('k_nearest_neigbour.csv', index = False)
```

**Output:**

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 135567 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 135568 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 135569 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 4 | 135570 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 135571 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 135572 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 135573 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 135574 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 135575 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 135576 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 135577 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 135578 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 13 | 135579 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 14 | 135580 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 15 | 135581 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 16 | 135582 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 135583 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 18 | 135584 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 19 | 135585 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 20 | 135586 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 21 | 135587 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 22 | 135588 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 23 | 135589 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 24 | 135590 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 25 | 135591 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 26 | 135592 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 27 | 135593 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 28 | 135594 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 29 | 135595 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 30 | 135596 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 31 | 135597 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 32 | 135598 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 33 | 135599 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 34 | 135600 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 35 | 135601 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 36 | 135602 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 37 | 135603 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 38 | 135604 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 39 | 135605 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 40 | 135606 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 41 | 135607 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 42 | 135608 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

k_nearest_neigbour

| 1902 | new | **NiteshJoshi** | 27.74876 | 1 |
|---|---|---|---|---|

**Your Best Entry ↑**
Congratulations on making your first submission!

For KNN we tried changing the K value ranging from 1 to 20. But we got the best results for these features when k=3.

## 2. Random Forest:

For this algorithm we have considered
- Category
- Dates
- x and y axis.

We took these features so that we might get the exact location, but when we tried it on real algorithm it gave us some astonishing results as follows.

Output:

| Id | ARSON | ASSAULT | BAD CHEC | BRIBERY | BURGLARY | DISORDER | DRIVING | DRUG/NA | DRUNKEN | EMBEZZLE | EXTORTIO | FAMILY OF | FORGERY/ | FRAUD | GAMBLIN( | KIDNAPPI | LARCENY/ | LIQUOR L | LOITERING | MISSI |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 3 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 6 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 11 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 12 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 17 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 21 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |

criminals

`

**Kaggle Ranking:**



1661　new　**AnuragPatsariya**　25.02448　1

**Your Best Entry ↑**
Congratulations on making your first submission!

# 3. Naïve bayes:
For the first run of kaggle submission we used Naïve Bayes with two feature.

**Feature used:**
Day of the week: As in our analyzed data we can see that crime rates become high on a specific day so we choose this feature for our algorithm

**PdDistrict:** Crime rates can be different based on the location, so instead of choosing exact X & Y coordinates for the location we choose PdDistrict so that crime location can be more generalized.

Process

Retrieve train2

inp

out

Naïve Bayes

tra
mod
exa

res

res

Set Role

exa
exa
ori

Apply Model

mod
lab
unl
mod

Retrieve test2

out

To leverage the Wisdom of Crowds, you must be a member of the RapidMiner Community!

After running the above model for Naïve Bayes algorithm we got the following output.

Calculating result: Distribution Table

Below is the small part of code to apply Naïve Bayes algorithm in python:

```
import pandas as pd
from sklearn.naive_bayes import BernoulliNB
.
.
features = ['Friday', 'Monday', 'Saturday', 'Sunday', 'Thursday', 'Tuesday',
 'Wednesday', 'BAYVIEW', 'CENTRAL', 'INGLESIDE', 'MISSION',
 'NORTHERN', 'PARK', 'RICHMOND', 'SOUTHERN', 'TARAVAL', 'TENDERLOIN']

model = BernoulliNB()
model.fit(train_data[features], train_data['crime'])
predicted = model.predict_proba(test_data[features])

#Write results

result=pd.DataFrame(predicted, columns=le_crime.classes_)
result.to_csv('NaivetestResult14.csv', index = True, index_label = 'Id' )
```

# Output

*(spreadsheet screenshot of NaivetestResult containing numeric probability values across crime category columns)*

Kaggle Ranking on submission:

**Vidhi Patel**

Verify account ›

NOVICE ✦ ?

Joined 16 minutes ago

Profile | Results | Scripts | Forum | Account | Activity

**San Francisco Crime Classification**
1 entry in team AnonymouS Team

Current
**1114th**/1990
Ending 27 days from now

**After getting the above results we tried the same code for three features given below:**

- **hour**
- **day of week**
- **PdDistrict**

For the second run, we used time, days of Week and PdDistrict.
This time we choose "Time" as a feature because there are possibilities of happening certain crimes at a particular time. e.g. drunk and drive happens specially in the night.

We changed the python code for Naïve Bayes as shown below:

```
import pandas as pd
from sklearn.naive_bayes import BernoulliNB
.
.
features = ['Friday', 'Monday', 'Saturday', 'Sunday', 'Thursday', 'Tuesday',
 'Wednesday', 'BAYVIEW', 'CENTRAL', 'INGLESIDE', 'MISSION',
 'NORTHERN', 'PARK', 'RICHMOND', 'SOUTHERN', 'TARAVAL', 'TENDERLOIN']

features2 = [x for x in range(0,24)]
features = features + features2

model = BernoulliNB()
model.fit(train_data[features], train_data['crime'])
predicted = model.predict_proba(test_data[features])

#Write results

result=pd.DataFrame(predicted, columns=le_crime.classes_)
result.to_csv('NaivetestResult14.csv', index = True, index_label = 'Id' )
```

Here we created sparse matrix of days of week, time and PdDistrict.

Output from algorithm:

Kaggle Ranking Run#2:

# 4.Logistic Regression

For this algorithm we used same feature which we selected for the second run of Naïve Bayes.

- Days of week
- PdDistrict
- Time

Output:

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 566655 | 0.00240019 | 0.07628601 | 0.00105801 | 0.00095714 | 0.0459254 | 0.00242661 | 0.00326742 | 0.0230667 | 0.00423561 | 0.00146854 | 0.00102723 | 0.00111286 | 0.00932607 | 0.01592356 | 0.0080496 | 0.00303885 | 0.21412268 | 0.00215246 | 0.00140473 | 0.04163648 | 0.083217 |
| 566656 | 0.00274129 | 0.06395558 | 0.00118143 | 0.00102055 | 0.05123252 | 0.00256674 | 0.00576235 | 0.02228219 | 0.00395507 | 0.0013804 | 0.00107916 | 0.00119412 | 0.00771104 | 0.01660481 | 0.00087478 | 0.00280153 | 0.25342881 | 0.00224009 | 0.00142635 | 0.021628 | 0.098966 |
| 566657 | 0.00425492 | 0.10309199 | 0.00092664 | 0.00105337 | 0.03884162 | 0.00228735 | 0.00252089 | 0.04795856 | 0.00329415 | 0.00128749 | 0.0008043 | 0.00117935 | 0.00539085 | 0.00739394 | 0.00082153 | 0.00340492 | 0.13931618 | 0.00174661 | 0.00130772 | 0.04607046 | 0.054660 |
| 566658 | 0.00274129 | 0.06395558 | 0.00118143 | 0.00102055 | 0.05123252 | 0.00256674 | 0.00576235 | 0.02228219 | 0.00395507 | 0.0013804 | 0.00107916 | 0.00119412 | 0.00771104 | 0.01660481 | 0.00087478 | 0.00280153 | 0.25342881 | 0.00224009 | 0.00142635 | 0.021628 | 0.098966 |
| 566659 | 0.00171765 | 0.06723915 | 0.00108612 | 0.00083182 | 0.04943369 | 0.00403249 | 0.00203962 | 0.02480294 | 0.00525614 | 0.00180066 | 0.00097829 | 0.00096888 | 0.00763113 | 0.01843731 | 0.00082372 | 0.00243252 | 0.35496388 | 0.00244651 | 0.00142125 | 0.01310758 | 0.099993 |
| 566660 | 0.0020035 | 0.061625 | 0.00106275 | 0.00101133 | 0.05431397 | 0.00401839 | 0.00336312 | 0.05753715 | 0.00676565 | 0.00148086 | 0.00096474 | 0.00123245 | 0.00636772 | 0.01378678 | 0.00084893 | 0.00252749 | 0.23425865 | 0.00363398 | 0.00150559 | 0.04196284 | 0.093808 |
| 566661 | 0.0020035 | 0.061625 | 0.00106275 | 0.00101133 | 0.05431397 | 0.00401839 | 0.00336312 | 0.05753715 | 0.00676565 | 0.00148086 | 0.00096474 | 0.00123245 | 0.00636772 | 0.01378678 | 0.00084893 | 0.00252749 | 0.23425865 | 0.00363398 | 0.00150559 | 0.04196284 | 0.093808 |
| 566662 | 0.0020035 | 0.061625 | 0.00106275 | 0.00101133 | 0.05431397 | 0.00401839 | 0.00336312 | 0.05753715 | 0.00676565 | 0.00148086 | 0.00096474 | 0.00123245 | 0.00636772 | 0.01378678 | 0.00084893 | 0.00252749 | 0.23425865 | 0.00363398 | 0.00150559 | 0.04196284 | 0.093808 |
| 566663 | 0.00245999 | 0.06008445 | 0.0011962 | 0.00103506 | 0.05507636 | 0.00255841 | 0.00527151 | 0.02642345 | 0.00389592 | 0.00145985 | 0.00108215 | 0.00123606 | 0.00768151 | 0.01666655 | 0.00087917 | 0.00272934 | 0.27062622 | 0.00286717 | 0.00141846 | 0.02065294 | 0.098265 |
| 566664 | 0.00383125 | 0.09743122 | 0.00094126 | 0.0010718 | 0.04193244 | 0.00228729 | 0.00231288 | 0.0567994 | 0.00325532 | 0.00136601 | 0.00080914 | 0.00122472 | 0.00538741 | 0.00744509 | 0.00082832 | 0.00332797 | 0.15073968 | 0.00224313 | 0.0013047 | 0.04419312 | 0.054403 |
| 566665 | 0.00168181 | 0.06485724 | 0.00091447 | 0.00077404 | 0.05080138 | 0.00308563 | 0.0229125 | 0.04739077 | 0.0034202 | 0.00132086 | 0.00076038 | 0.0007975 | 0.00708121 | 0.01336505 | 0.00066052 | 0.00226656 | 0.32779486 | 0.00198675 | 0.00191665 | 0.01488827 | 0.075356 |
| 566666 | 0.00171765 | 0.06723915 | 0.00108612 | 0.00083182 | 0.04943369 | 0.00403249 | 0.00203962 | 0.02480294 | 0.00525614 | 0.00180066 | 0.00097829 | 0.00096888 | 0.00763113 | 0.01843731 | 0.00082372 | 0.00243252 | 0.35496388 | 0.00244651 | 0.00142125 | 0.01310758 | 0.099993 |
| 566667 | 0.00168181 | 0.06485724 | 0.00091447 | 0.00077404 | 0.05080138 | 0.00308563 | 0.0229125 | 0.04739077 | 0.0034202 | 0.00132086 | 0.00076038 | 0.0007975 | 0.00708121 | 0.01336505 | 0.00066052 | 0.00226656 | 0.32779486 | 0.00198675 | 0.00191665 | 0.01488827 | 0.075356 |
| 566668 | 0.00168181 | 0.06485724 | 0.00091447 | 0.00077404 | 0.05080138 | 0.00308563 | 0.0229125 | 0.04739077 | 0.0034202 | 0.00132086 | 0.00076038 | 0.0007975 | 0.00708121 | 0.01336505 | 0.00066052 | 0.00226656 | 0.32779486 | 0.00198675 | 0.00191665 | 0.01488827 | 0.075356 |
| 566669 | 0.00171765 | 0.06723915 | 0.00108612 | 0.00083182 | 0.04943369 | 0.00403249 | 0.00203962 | 0.02480294 | 0.00525614 | 0.00180066 | 0.00097829 | 0.00096888 | 0.00763113 | 0.01843731 | 0.00082372 | 0.00243252 | 0.35496388 | 0.00244651 | 0.00142125 | 0.01310758 | 0.099993 |
| 566670 | 0.00148259 | 0.06635384 | 0.00081998 | 0.00071679 | 0.02934291 | 0.00251801 | 0.00199237 | 0.06476474 | 0.00549706 | 0.00154388 | 0.00068591 | 0.00076629 | 0.00844922 | 0.01501734 | 0.00058565 | 0.0022883 | 0.328808 | 0.00290218 | 0.00256508 | 0.01552875 | 0.099053 |
| 566671 | 0.00148259 | 0.06635384 | 0.00081998 | 0.00071679 | 0.02934291 | 0.00251801 | 0.00199237 | 0.06476474 | 0.00549706 | 0.00154388 | 0.00068591 | 0.00076629 | 0.00844922 | 0.01501734 | 0.00058565 | 0.0022883 | 0.328808 | 0.00290218 | 0.00256508 | 0.01552875 | 0.099053 |
| 566672 | 0.00148259 | 0.06635384 | 0.00081998 | 0.00071679 | 0.02934291 | 0.00251801 | 0.00199237 | 0.06476474 | 0.00549706 | 0.00154388 | 0.00068591 | 0.00076629 | 0.00844922 | 0.01501734 | 0.00058565 | 0.0022883 | 0.328808 | 0.00290218 | 0.00256508 | 0.01552875 | 0.099053 |
| 566673 | 0.00171765 | 0.06723915 | 0.00108612 | 0.00083182 | 0.04943369 | 0.00403249 | 0.00203962 | 0.02480294 | 0.00525614 | 0.00180066 | 0.00097829 | 0.00096888 | 0.00763113 | 0.01843731 | 0.00082372 | 0.00243252 | 0.35496388 | 0.00244651 | 0.00142125 | 0.01310758 | 0.099993 |
| 566674 | 0.00245999 | 0.06008445 | 0.0011962 | 0.00103506 | 0.05507636 | 0.00255841 | 0.00527151 | 0.02642345 | 0.00389592 | 0.00145985 | 0.00108215 | 0.00123606 | 0.00768151 | 0.01666655 | 0.00087917 | 0.00272934 | 0.27062622 | 0.00286717 | 0.00141846 | 0.02065294 | 0.098265 |
| 566675 | 0.00148259 | 0.06635384 | 0.00081998 | 0.00071679 | 0.02934291 | 0.00251801 | 0.00199237 | 0.06476474 | 0.00549706 | 0.00154388 | 0.00068591 | 0.00076629 | 0.00844922 | 0.01501734 | 0.00058565 | 0.0022883 | 0.328808 | 0.00290218 | 0.00256508 | 0.01552875 | 0.099053 |
| 566676 | 0.00171765 | 0.06723915 | 0.00108612 | 0.00083182 | 0.04943369 | 0.00403249 | 0.00203962 | 0.02480294 | 0.00525614 | 0.00180066 | 0.00097829 | 0.00096888 | 0.00763113 | 0.01843731 | 0.00082372 | 0.00243252 | 0.35496388 | 0.00244651 | 0.00142125 | 0.01310758 | 0.099993 |
| 566677 | 0.00171765 | 0.06723915 | 0.00108612 | 0.00083182 | 0.04943369 | 0.00403249 | 0.00203962 | 0.02480294 | 0.00525614 | 0.00180066 | 0.00097829 | 0.00096888 | 0.00763113 | 0.01843731 | 0.00082372 | 0.00243252 | 0.35496388 | 0.00244651 | 0.00142125 | 0.01310758 | 0.099993 |
| 566678 | 0.00239647 | 0.09410002 | 0.00097796 | 0.00111076 | 0.04028616 | 0.00220043 | 0.00265458 | 0.03487087 | 0.00305875 | 0.00132603 | 0.0009263 | 0.00128539 | 0.00723754 | 0.0109684 | 0.00080261 | 0.00391274 | 0.17040947 | 0.00234182 | 0.00123859 | 0.03193286 | 0.069266 |
| 566679 | 0.00216312 | 0.05496454 | 0.00114943 | 0.00096877 | 0.06455522 | 0.00267768 | 0.00435731 | 0.02734554 | 0.00342822 | 0.00149011 | 0.00103475 | 0.00115374 | 0.0089095 | 0.01847745 | 0.00086101 | 0.00274804 | 0.25985257 | 0.00284422 | 0.00148017 | 0.02221482 | 0.100173 |
| 566680 | 0.0013113 | 0.06109011 | 0.00079261 | 0.00067487 | 0.03477468 | 0.00265114 | 0.0016557 | 0.06732119 | 0.0048669 | 0.00158529 | 0.00065977 | 0.00071951 | 0.00985739 | 0.01675221 | 0.00057698 | 0.00223726 | 0.31846484 | 0.00289613 | 0.00269249 | 0.01681045 | 0.101581 |
| 566681 | 0.00211322 | 0.08657096 | 0.00094237 | 0.00104258 | 0.04749934 | 0.00230961 | 0.0021994 | 0.03617897 | 0.00269889 | 0.00135736 | 0.00088823 | 0.0012032 | 0.00841932 | 0.01220337 | 0.00078827 | 0.0039507 | 0.16320875 | 0.00232966 | 0.00129615 | 0.0344151 | 0.070870 |
| 566682 | 0.00211322 | 0.08657096 | 0.00094237 | 0.00104258 | 0.04749934 | 0.00230961 | 0.0021994 | 0.03617897 | 0.00269889 | 0.00135736 | 0.00088823 | 0.0012032 | 0.00841932 | 0.01220337 | 0.00078827 | 0.0039507 | 0.16320875 | 0.00232966 | 0.00129615 | 0.0344151 | 0.070870 |
| 566683 | 0.00152097 | 0.06198687 | 0.00105107 | 0.00078407 | 0.05841768 | 0.00425019 | 0.00169692 | 0.02585286 | 0.00465882 | 0.00185105 | 0.00094209 | 0.00091078 | 0.00891413 | 0.02058174 | 0.00081245 | 0.00246663 | 0.34475421 | 0.00244419 | 0.00149363 | 0.0142082 | 0.102655 |
| 566684 | 0.00383503 | 0.0895497 | 0.00090596 | 0.00100483 | 0.04937106 | 0.00239796 | 0.00191393 | 0.05881065 | 0.00286904 | 0.00139666 | 0.0007498 | 0.00114506 | 0.00626175 | 0.00827712 | 0.00081257 | 0.00335635 | 0.14404959 | 0.00222888 | 0.00136374 | 0.04752667 | 0.055621 |
| 566685 | 0.00211322 | 0.08657096 | 0.00094237 | 0.00104258 | 0.04749934 | 0.00230961 | 0.0021994 | 0.03617897 | 0.00269889 | 0.00135736 | 0.00088823 | 0.0012032 | 0.00841932 | 0.01220337 | 0.00078827 | 0.0039507 | 0.16320875 | 0.00232966 | 0.00129615 | 0.0344151 | 0.070870 |
| 566686 | 0.00383503 | 0.0895497 | 0.00090596 | 0.00100483 | 0.04937106 | 0.00239796 | 0.00191393 | 0.05881065 | 0.00286904 | 0.00139666 | 0.0007498 | 0.00114506 | 0.00626175 | 0.00827712 | 0.00081257 | 0.00335635 | 0.14404959 | 0.00222888 | 0.00136374 | 0.04752667 | 0.055621 |
| 566687 | 0.00211322 | 0.08657096 | 0.00094237 | 0.00104258 | 0.04749934 | 0.00230961 | 0.0021994 | 0.03617897 | 0.00269889 | 0.00135736 | 0.00088823 | 0.0012032 | 0.00841932 | 0.01220337 | 0.00078827 | 0.0039507 | 0.16320875 | 0.00232966 | 0.00129615 | 0.0344151 | 0.070870 |
| 566688 | 0.00383503 | 0.0895497 | 0.00090596 | 0.00100483 | 0.04937106 | 0.00239796 | 0.00191393 | 0.05881065 | 0.00286904 | 0.00139666 | 0.0007498 | 0.00114506 | 0.00626175 | 0.00827712 | 0.00081257 | 0.00335635 | 0.14404959 | 0.00222888 | 0.00136374 | 0.04752667 | 0.055621 |
| 566689 | 0.00149081 | 0.05984517 | 0.00088589 | 0.00073037 | 0.06007445 | 0.0032558 | 0.00190838 | 0.04940158 | 0.00303405 | 0.00135927 | 0.00073301 | 0.00075046 | 0.00828112 | 0.01494395 | 0.00065217 | 0.00230076 | 0.31831615 | 0.00198695 | 0.00201632 | 0.01615275 | 0.07749 |
| 566690 | 0.00149081 | 0.05984517 | 0.00088589 | 0.00073037 | 0.06007445 | 0.0032558 | 0.00190838 | 0.04940158 | 0.00303405 | 0.00135927 | 0.00073301 | 0.00075046 | 0.00828112 | 0.01494395 | 0.00065217 | 0.00230076 | 0.31831615 | 0.00198695 | 0.00201632 | 0.01615275 | 0.07749 |
| 566691 | 0.0013113 | 0.06109011 | 0.00079261 | 0.00067487 | 0.03477468 | 0.00265114 | 0.0016557 | 0.06732119 | 0.0048669 | 0.00158529 | 0.00065977 | 0.00071951 | 0.00985739 | 0.01675221 | 0.00057698 | 0.00223726 | 0.31846484 | 0.00289613 | 0.00269249 | 0.01681045 | 0.101581 |
| 566692 | 0.0013113 | 0.06109011 | 0.00079261 | 0.00067487 | 0.03477468 | 0.00265114 | 0.0016557 | 0.06732119 | 0.0048669 | 0.00158529 | 0.00065977 | 0.00071951 | 0.00985739 | 0.01675221 | 0.00057698 | 0.00223726 | 0.31846484 | 0.00289613 | 0.00269249 | 0.01681045 | 0.101581 |
| 566693 | 0.0017558 | 0.0561846 | 0.00101768 | 0.00094331 | 0.06346114 | 0.00419101 | 0.00276939 | 0.05926648 | 0.00593505 | 0.00150637 | 0.00091932 | 0.00114643 | 0.00736214 | 0.01523794 | 0.00082855 | 0.0025361 | 0.22370924 | 0.00359255 | 0.00156569 | 0.04490304 | 0.095316 |
| 566694 | 0.00152097 | 0.06198687 | 0.00105107 | 0.00078407 | 0.05841768 | 0.00425019 | 0.00169692 | 0.02585286 | 0.00465882 | 0.00185105 | 0.00094209 | 0.00091078 | 0.00891413 | 0.02058174 | 0.00081245 | 0.00246663 | 0.34475421 | 0.00244419 | 0.00149363 | 0.0142082 | 0.102655 |
| 566695 | 0.00138981 | 0.07271132 | 0.00081885 | 0.00089289 | 0.03543459 | 0.00637327 | 0.00219625 | 0.08514825 | 0.00506784 | 0.00125194 | 0.00071593 | 0.00126503 | 0.00742287 | 0.01297071 | 0.00066183 | 0.00259828 | 0.18978838 | 0.00441107 | 0.00188162 | 0.02385428 | 0.08448 |

# 6. Conclusion:

For this competition we tried these four algorithms. Out of which logistic regression with three features

- PdDistrict
- Hour
- Days of week

Gives us the best results and out ranking bumped up to 875 from 1902.

| Algorithm | Kaggle Ranking |
|---|---|
| K Nearest Neighbor (k = 3) | 1902 |
| Random Forest | 1667 |
| Naïve Bayes | 1114 |
| Naïve Bayes (three features) | 910 |
| Logistic Regression | 875 |

We can conclude that good feature selection and good algorithm can improve the ranking at Kaggle competition.

# 7. References:

1. http://efavdb.com/predicting-san-francisco-crimes/
2. http://scikit-learn.org/stable/
3. https://www.kaggle.com/c/sf-crime
4. http://docs.rapidminer.com/studio/operators/modeling/predictive/bayesian/naive_bayes.html