# Project Poirot: Machine Learning based Predictive Crime Analysis
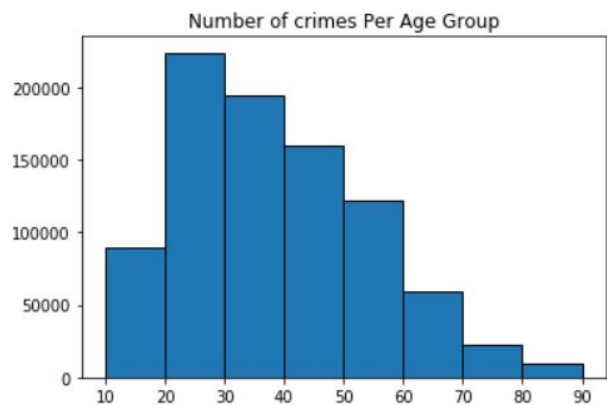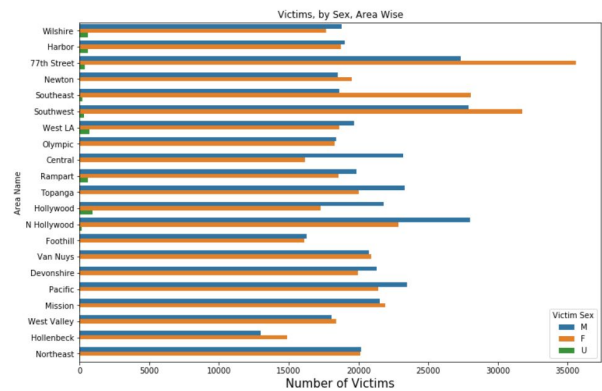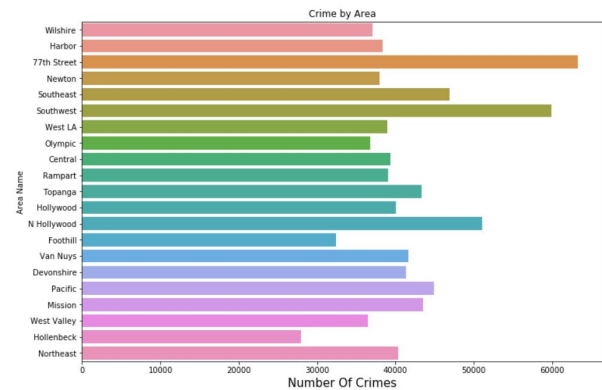
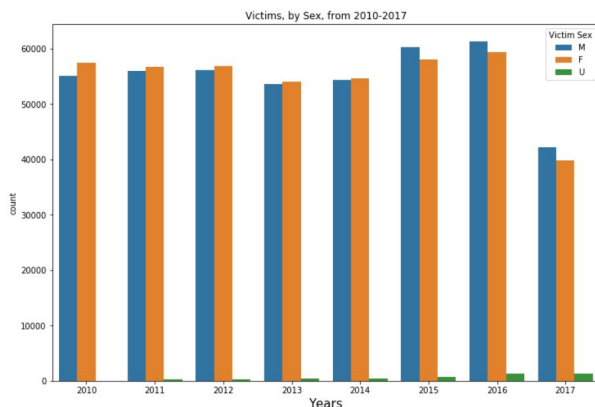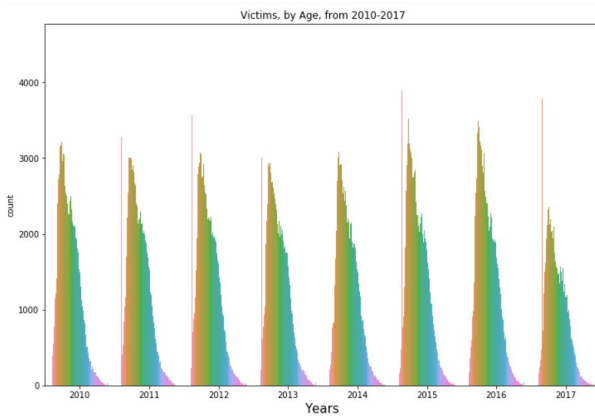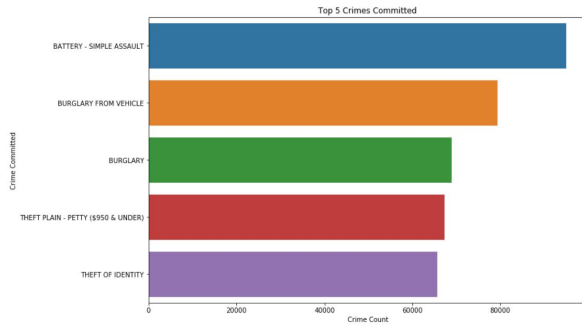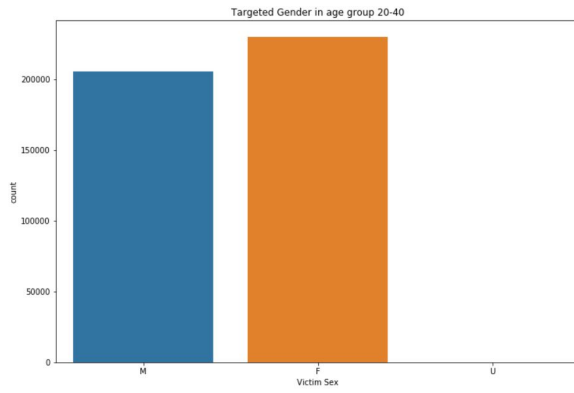Mansi Ganatra[1], Naga Ritwik Indugu [2] and Uma Kanumuri [3]

*Abstract*—LA is a very beautiful city, but one thing that scares a lot of people is the amount of crime that happens here. Our intention is to understand what makes this beautiful city to be infested with crime.To work upon this we have collected the dataset from LAPD official site.The dataset contains approximately 1.8 Million rows with around 28 attributes including DR Number, Date Reported, Date Occurred, Time Occurred, Area ID, Area Name Reporting District etc.Some of the questions which we are going to answer are Predicting the time of day, the days of a week, the days of a month, the days of a year in which different kinds of crimes are most likely to happen,to predict the hot spots for various crime categories, to predict the crime patterns, To predict the probability of an offender repeating the crime,Crimes correlation with age,Severity of crime by area (which crime is highest in each area) etc using the different machine learning algorithms like Decision Trees,Nave Bayes,Logistic Regression,K-Nearest Neighbours(K-NN),Support Vector Machines(SVM),Clustering Algorithms.The libraries which we are using include pandas, numpy, scikit-learn, pytorch, theano, etc.For the visualization of results the libraries being used are matplotlib, d3.js etc.

## I. INTRODUCTION

Beautiful and historic, Los Angeles is also infested with a high amount of crime. The total crimes reported to LAPD in 2017 were 129587(Violent: 29661, Property: 99926). The LAPD has made available the crime dataset from 2010-2018. In the Machine Learning experiments described below, we use the same crime dataset pre-processed from Kaggle containing records from 2010-2017. We performed some exploratory data analysis on the features and the results are as in the following section. In this project, we try to answer different questions based on the features. The questions and the algorithms used to evaluate answers to each question are described in the further section. We have used accuracy, F-score, mean squared error as our performance metrics. There have been multiple questions raised for the ethicality of such studies. The study in this project is only empirical and academic in nature. We do not intend to use the results and predictions on field or in conjunction with other predictive policing technologies in use today. Our aim is to apply various machine learning techniques to the selected features, keeping into account the literature survey, and evaluate answer to the selected questions. We intend to measure the accuracy of our results using the performance metric as described in the following section.

## II. EXPLORATORY DATA ANALYSIS

Targeted Gender in age group 20-40



Top 5 Crimes Committed



Victims, by Age, from 2010-2017



Victims, by Sex, from 2010-2017

## III. METHODOLOGY

### A. Learning Algorithms

The template is used to format your paper and style the text. All margins, column widths, line spaces, and text fonts are prescribed; please do not alter them. You may note peculiarities. For example, the head margin in this template measures proportionately more than is customary. This measurement and others are deliberate, using specifications that anticipate your paper as one part of the entire proceedings, and not as an independent document. Please do not revise any of the current designations

- **SVM**: In this project we have used SVM with rbf kernel keeping the value of c to 1.0 and gamma to auto. Predict(X) method is used in predicting the crime code given victim age,victim sex.sklearn.svm has been used in importing the svm function.
- **Decision Trees**:Decision trees with a criteria of entropy and random state as 0 are taken.sklearn.tree.decision tree classifier is used to get the decision tree function.Predict(X) method has been used in predicting the crime code given victim age and victim sex
- **Logistic Regression**:Mutinomial logistic regression has been used in predicting the crime code given victim age and victim sex.Newton Solver has been used here. sklearn.linearmodel.logistic regression has been used for getting the regression function.Predict(X) method has been used for predicting the crime code.
- **Linear Regression**:Linear regression has been used in predicting the time occured given the weekdays.Here weekday is calculated from the date occurred using a python function.sklearn.linearmodel.linearregression has been used for getting the prediction function. predict(X) method has been used for predicting the linear regression.
- **Neural Networks**:Neural network with 3 layers and approxiamately 10000 weights have been used in predicting the time occurred given weekday and crimecode.Also neural networks have been used in predicting the crimecode given victim age and gender.

### B. Performance Metrics

- **Accuracy score**:From sklearn.metrics import accuracy score has been used in getting the accuracy score function.It is used for calculating the accuracy's of svm, decision trees and logistic regression in this project.
- **Root mean square error**:RMSE has been used in linear regression for measuring the accuracy of predicting time of the day given day and time. from sklearn.metrics import mean_squared_error has been used for getting the mean_squared_error function.

## IV. VISUALIZATION FUNCTIONS

Packages seaborn and matplotlib in python have been used in visualization of the results.

## A. Prediction of time occurred given weekday and Crime Code

- **Preprocessing**:Weekday has been obtained from the date of the crime occurred with python function dt.dayofweek.
- **Linear regression**:Linear regression has been used for the prediction given only the weekday. Accuracy obtained was:

(1).png (1).png



- **Logistic Regression**:

(3).png (3).png



- **Decision Trees**:

(4).png (4).png



- **K-NN**:

(5).png (5).png



- **Neural Networks**: A neural network of sequential model with 3 layers and approximately 18000 weights have been built with relu,tanh and softmax functions. Accuracy obtained was:



We later introduced a new input feature specifying whether the day was weekend or not. Accuracy obtained was:

| Name | Type | Size | Value |
|---|---|---|---|
| X | int64 | (1048575, 3) | Min: 0<br>Max: 956 |
| X_test | int64 | (209715, 3) | Min: 0<br>Max: 956 |
| X_train | int64 | (838860, 3) | Min: 0<br>Max: 956 |
| Y | object | (1048575, 1) | Min: '00'<br>Max: '23' |
| accuracies | list | 1 | [7.996089931577355] |
| b | DataFrame | (1048575, 5) | Column names: Date Occurred, Time Occurred, Crime Code, weekday, month |
| column_1 | Series | (1048575,) | Series object of pandas.core.series module |
| | DataFrame | (1048575, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |

```
dense_3 (Dense)              (None, 113)        21470
dense_4 (Dense)              (None, 24)         2736
=================================================================
Total params: 72,645
Trainable params: 72,645
Non-trainable params: 0

In [3]: X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2,
random_state=0)

In [4]: model.fit(X_train, y_train, epochs=150, batch_size=10, verbose=0)
   ...: scores = model.evaluate(X_test, y_test, verbose=0)
   ...: print("%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))
   ...: accuracies.append(scores[1] * 100)
   ...: print("%.2f%% (+/- %.2f%%)" % (np.mean(accuracies), np.std(accuracies)))
acc: 8.00%
8.00% (+/- 0.00%)

In [5]:
```

After this we introduced another input feature specifying the season corresponding to the day of crime. Accuracy obtained was:

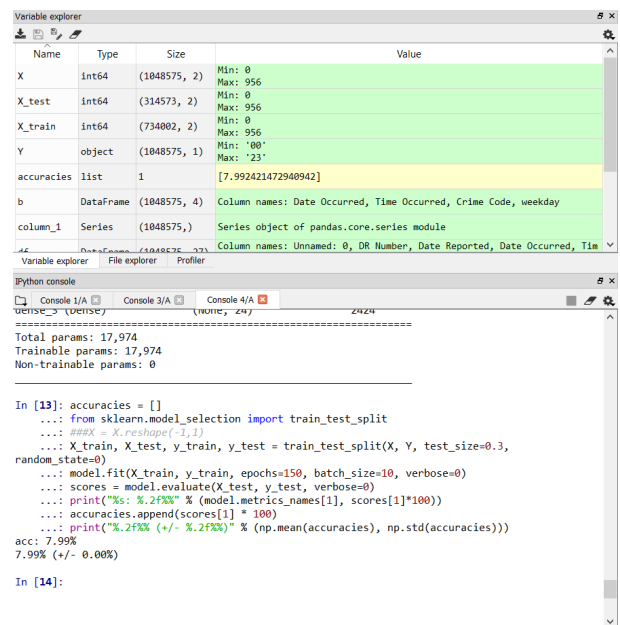| Name | Type | Size | Value |
|---|---|---|---|
| X | int64 | (1048575, 3) | Min: 0<br>Max: 956 |
| X_test | int64 | (209715, 3) | Min: 0<br>Max: 956 |
| X_train | int64 | (838860, 3) | Min: 0<br>Max: 956 |
| Y | object | (1048575, 1) | Min: '00'<br>Max: '23' |
| accuracies | list | 1 | [5.9218463151805025] |
| b | DataFrame | (1048575, 5) | Column names: Date Occurred, Time Occurred, Crime Code, weekday, month |
| column_1 | Series | (1048575,) | Series object of pandas.core.series module |
| | DataFrame | (1048575, 27) | Column names: Unnamed: 0, DR Number, Date Occurred, Tim ... |

```
In [34]: from sklearn.model_selection import train_test_split
   ...: ###X = X.reshape(-1,1)
   ...: X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2,
random_state=0)

In [35]: model.fit(X_train, y_train, epochs=150, batch_size=10, verbose=0)
   ...: scores = model.evaluate(X_test, y_test, verbose=0)
   ...: print("%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))
   ...: accuracies.append(scores[1] * 100)
   ...: print("%.2f%% (+/- %.2f%%)" % (np.mean(accuracies), np.std(accuracies)))
acc: 5.92%
5.92% (+/- 0.00%)

In [36]:
```

Permissions: RW    End-of-lines: CRLF    Encoding: ASCII    Line: 42    Column: 17    Memory: 52 %

## B. Prediction of Crime Code given victim information

- **Preprocessing**: The data had Victim Age range from 0-99. For better processing, we clustered the age into groups of 10.
- **Logistic Regression**:

| Name | Type | Size | Value |
|---|---|---|---|
| X | float64 | (328775, 3) | [[1. 0. 4.]<br>[1. 0. 2.]] |
| X_test | float64 | (82194, 3) | [[ 1.0719946  -0.05426602  1.31268955]<br>[-0.93284052 -0.05426602 -0.67 ... |
| X_train | float64 | (246581, 3) | [[-0.93284052 -0.05426602 -0.67696095]<br>[ 1.0719946  -0.05426602 -0.67 ... |
| ac | float64 | 1 | 0.4920432148331995 |
| dataset | DataFrame | (328775, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |
| y | int64 | (328775,) | [6 2 2 ... 6 2 6] |
| y_pred | int64 | (82194,) | [2 6 6 ... 2 2 2] |

- **Decision Trees**:

| Name | Type | Size | Value |
|---|---|---|---|
| X | float64 | (328775, 4) | [[4. 0. 1. 0.]<br>[2. 0. 1. 0.]] |
| X_test | float64 | (82194, 4) | [[ 1.31268955 -1.06569017  1.0719946  -0.05426602]<br>[-0.67696095  0.93 ... |
| X_train | float64 | (246581, 4) | [[-0.67696095  0.93835903 -0.93284052 -0.05426602]<br>[-0.67696095 -1.06 ... |
| ac | float64 | 1 | 0.4943061537338492 |
| dataset | DataFrame | (328775, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |
| train | DataFrame | (328775, 5) | Column names: Victim Age, Crime Code, Gender_F, Gender_M, Gender_U |
| y | int64 | (328775,) | [6 2 2 ... 6 2 6] |

- **Neural Networks**: A neural network of sequential model with 3 layers and approximately 18000 weights have been built with relu,tanh and softmax functions. Accuracy obtained was:

| Name | Type | Size | Value |
|---|---|---|---|
| X_full | float64 | (880758, 5) | Min: 0.0<br>Max: 21.0 |
| X_test | float64 | (220190, 5) | Min: 0.0<br>Max: 21.0 |
| X_train | float64 | (660568, 5) | Min: 0.0<br>Max: 21.0 |
| accuracies | list | 2 | [34.3253079168231, 34.40256142427867] |
| dataset | DataFrame | (880758, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |
| scores | list | 2 | [1.7368069640318604, 0.3440256142427867] |
| train | DataFrame | (880758, 6) | Column names: Crime Code, Victim Age, Area ID, Gender_F, Gender_M, Gen ... |
| | int64 | (880758,) | Min: 1 |

```
660568/660568 [==============================] - 288s 436us/step - loss: 1.7403 - acc: 0.3404
Epoch 4/10
660568/660568 [==============================] - 256s 388us/step - loss: 1.7435 - acc:
0.34177990/660568 [========================>......] - ETA: 43s - loss: 1.7436 - acc: 0.3415
Epoch 5/10
660568/660568 [==============================] - 263s 398us/step - loss: 1.7419 - acc: 0.3423
Epoch 6/10
660568/660568 [==============================] - 260s 394us/step - loss: 1.7400 - acc: 0.3428
ETA: 1:48 - loss: 1.7403 - acc: 0.3428 - ETA: 1:23 - loss: 1.7400 - acc: 0.3428
Epoch 7/10
660568/660568 [==============================] - 250s 378us/step - loss: 1.7391 - acc: 0.3430
Epoch 8/10
660568/660568 [==============================] - 269s 407us/step - loss: 1.7382 - acc: 0.3436
Epoch 9/10
660568/660568 [==============================] - 258s 390us/step - loss: 1.7377 - acc: 0.3432
Epoch 10/10
660568/660568 [==============================] - 251s 380us/step - loss: 1.7394 - acc: 0.3430
Evaluating train dataset:
acc: 34.33%
Evaluating dev dataset:
acc: 34.40%
34.36% (+/- 0.04%)

In [9]:
```

**Variable explorer**

| Name | Type | Size | Value |
|---|---|---|---|
| X_full | float64 | (880758, 25) | Min: 0.0 Max: 21.0 |
| X_test | float64 | (220190, 25) | Min: 0.0 Max: 21.0 |
| X_train | float64 | (660568, 25) | Min: 0.0 Max: 21.0 |
| accuracies | list | 2 | [35.86337818369598, 36.0179844681219] |
| dataset | DataFrame | (880758, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |
| scores | list | 2 | [1.7183226979868662, 1.36017984468121905] |
| train | DataFrame | (880758, 26) | Column names: Crime Code, Victim Age, Area ID, Gender_F, Gender_M, Gen ... |
| | | | Min: 1 |

Variable explorer    File explorer    Profiler

**IPython console**

Console 6/A    Console 4/A    Console 5/A    Console 7/A    Console 8/A    Console 9/A

```
                                          - 310s 478us/step - loss: 1.7238 - acc: 0.3304
ETA: 48s - loss: 1.7264 - acc: 0.3563
Epoch 4/10
660568/660568 [==============================] - 246s 372us/step - loss: 1.7228 - acc: 0.3573
Epoch 5/10
660568/660568 [==============================] - 252s 381us/step - loss: 1.7216 - acc: 0.3578
Epoch 6/10
660568/660568 [==============================] - 256s 388us/step - loss: 1.7207 - acc: 0.3583
ETA: 1:39 - loss: 1.7197 - acc: 0.3592
Epoch 7/10
660568/660568 [==============================] - 248s 375us/step - loss: 1.7203 - acc: 0.3581
Epoch 8/10
660568/660568 [==============================] - 245s 371us/step - loss: 1.7200 - acc: 0.3582
Epoch 9/10
660568/660568 [==============================] - 265s 402us/step - loss: 1.7197 - acc: 0.3580
Epoch 10/10
660568/660568 [==============================] - 261s 395us/step - loss: 1.7196 - acc: 0.3580
Evaluating train dataset:
acc: 35.86%
Evaluating dev dataset:
acc: 36.02%
35.94% (+/- 0.08%)

In [15]:
```

### C. Predicting Crime Categories

We categorized crimes into two major subtypes: violent and non-violent. We tried to predict the crime category using the victim information. This can provide useful insights into what type of crime each age group is targeted to. Below are the results:

- **Decision Trees**:

| ac | float64 | 1 | 77.43819349815068 |
|---|---|---|---|

- **K-NN**:

| ac | float64 | 1 | 71.37312633832977 |
|---|---|---|---|

- **Neural Networks**:

| accuracies | list | 2 | [77.41926326089961, 77.417510121997275] |
|---|---|---|---|

### D. Using the 'Weapon Used Code' feature

The data contains information regarding the weapon used for crime: Weapon Used Code, it's description. EDA indicated that Weapon Used Code is highly correlated to the type of crime committed. But, it is fairly easy to deduce the crime type if we know the type of weapon used in it. Below results prove that point, as accuracies to the previous questions drastically shoot up once this information is introduced in the input features.

- **Prediction of Crime Code given victim information**
  - **Logistic Regression**:

**Variable explorer**

| Name | Type | Size | Value |
|---|---|---|---|
| X | float64 | (328775, 4) | [[1. 0. 4. 4.] [1. 0. 2. 1.] |
| X_test | float64 | (82194, 4) | [[ 1.0719946  -0.05426602  1.31268955   0.26908381] [-0.93284052 -0.05 ... |
| X_train | float64 | (246581, 4) | [[-0.93284052 -0.05426602 -0.67696095  1.14897659] [ 1.0719946  -0.05 ... |
| ac | float64 | 1 | 0.6064554590359393 |
| dataset | DataFrame | (328775, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |
| y | int64 | (328775,) | [6 2 2 ... 6 2 6] |
| y_pred | int64 | (82194,) | [6 6 6 ... 6 6 6] |

Variable explorer    File explorer    Help

- **Prediction of time occurred given weekday and Crime Code**
  - **Decision Trees**:

**Variable explorer**

| Name | Type | Size | Value |
|---|---|---|---|
| classification_reports | list | 0 | [] |
| column_1 | Series | (347265L,) | Series object of pandas.core.series module |
| confusion_matrices | list | 0 | [] |
| df | DataFrame | (347265, 27) | Column names: Unnamed: 0, DR Number, Date Reported, Date Occurred, Tim ... |
| seed | int | 1 | 5 |
| y_test | uint8 | (69453L,) | Min: 0 Max: 1 |
| y_train | uint8 | (277812L,) | Min: 0 Max: 1 |

**IPython console**

Console 8/A    Console 10/A    Console 11/A

```
  b['month'] = column_1.dt.month
D:/Sem1/INF552/Project/code/crime_weekday_season_time_decisiontrees.py:97:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/
indexing.html#indexing-view-versus-copy
  b['isWeekend'] = b['weekday'].apply(isWeekday)
D:/Sem1/INF552/Project/code/crime_weekday_season_time_decisiontrees.py:98:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/
indexing.html#indexing-view-versus-copy
  b['season']=b['month'].apply(season_crime)
C:\ProgramData\Anaconda2\lib\site-packages\sklearn\utils\validation.py:475:
DataConversionWarning: Data with input dtype uint8 was converted to float64 by
StandardScaler.
  warnings.warn(msg, DataConversionWarning)
*********************** DT Accuracy:
0.7537183419002779
<sklearn.tree._tree.Tree object at 0x0000000028148098>
In [2]:
```

## VI. RELEVANT WORK

The pioneer work in this field has been done by the team behind what is today known as PredPol. The introductory paper "Randomized Controlled Field Trials of Predictive Policing"[1] was first published in 2015. PredPol uses Crime type, Crime location, Crime date and time features and proprietary machine learning algorithms for forecasting. Predictive Policing: The Role of Crime Forecasting in Law Enforcement Operations [2] details the predictive policing strategies and their use in Law Enforcement. Predictive Policing: A Review of the Literature [3] from Portland State University reviews all the literature on predictive policing. Recent advancements include Partially Generative Neural Networks for Gang Crime Classification with Partial Information [4].

## VII. INDIVIDUAL CONTRIBUTIONS

- **Naga Ritwik Indugu**: Prediction of Crime Code given victim information
- **Uma Kanumuri**: Prediction of time occurred given weekday and Crime Code
- **Mansi Ganatra**: Predicting Crime Categories

## ACKNOWLEDGMENT

## REFERENCES

[1] G. O. Mohler, M. B. Short, Sean Malinowski, Mark Johnson, G. E. Tita, Andrea L. Bertozzi & P. J. Brantingham (2015) Randomized Controlled Field Trials of Predictive Policing, Journal of the American Statistical Association, 110:512, 1399-1411, DOI: 10.1080/01621459.2015.1077710

[2] Perry, Walter L., Brian McInnis, Carter C. Price, Susan Smith, and John S. Hollywood, Predictive Policing: The Role of Crime Forecasting in Law Enforcement Operations. Santa Monica, CA: RAND Corporation, 2013. https://www.rand.org/pubs/research_reports/RR233.html. Also available in print form.

[3] Portland State University. Criminology and Criminal Justice Senior Capstone, "Predictive Policing: A Review of the Literature" (2012). Criminology and Criminal Justice Senior Capstone Project. 5.

[4] Seo, Sungyong, Hau Kong Chan, P. Jeffrey Brantingham, Jorja Leap, Phebe Vayanos, Milind Tambe and Yan Liu. Partially Generative Neural Networks for Gang Crime Classification with Partial Information. (2017).

[5] http://michael-harmon.com/blog/crimetime.html

[6] https://machinelearningmastery.com/