# COURSERA DATASCIENCE CAPSTONE

Anshuman H

# INTRODUCTION

- Predicting Severity of Accident is important
  - To understand the factors leading to severity
  - Understand measure which can be implemeneted
  - Testing the impact of measures

# DATA ACQUISITION AND CLEANING

- Data file used is "Collision - All Years" which has been sourced from SDOT Traffic Management Division.

- The data has 38 attributes and covers the time period from 2004 to present date.

- The data is split into 2 kinds of severity "Injury collision" and "Property Damage Only Collision".
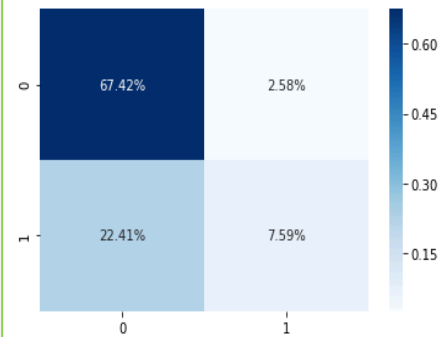
# METHODOLOGY

- Recursive Feature Elimination (RFE) has been used to further refine the selected features
  - Predictor_Matrix_N=['PERSONCOUNT','PEDCOUNT','PEDCYLCOUNT','VEHCOUNT','SDOT_COLCODE', 'INTKEY','LIGHTCOND_N','ROADCOND_N','WEATHER_N', 'INATTENTIONIND_N','UNDERINFL_N','HITPARKEDCAR_N','SPEEDING_N', 'COLLISIONTYPE_N','JUNCTIONTYPE_N','PEDROWNOTGRNT_N']
- The performance of the below models would be compared
  - Logistic Regression
  - Gaussian Naïve Bayes
  - Random Forest Classifier

# RESULT

# CONCLUSION

- Based on the dataset a model to predict severity was initially developed using LogisticRegression and then further modelled using the below 3 models.

    - Logistic Regression

    - Gaussian Naïve Bayes

    - Random Forest Classifier

- **Random Forest Classifier model was able to achieve 75% success rate based on the test data.**