# Executive Summary

I.  Data Collection and Cleaning

The Washington Post gathered data on criminal homicides from 50 major American cities over a decade, detailing incident locations, arrest status, and victim demographics. The dataset underwent meticulous cleaning and standardization, involving months of work to ensure accuracy. The primary dataset, obtained from Kaggle, encompassed information on over 52,000 homicides. Additionally, predicted population data from the US Census was incorporated to analyze the relationship between city size and homicide rates.

II.  Machine Learning and Model Performance

We employed the LGBMClassifier algorithm to create a predictive model. The model demonstrated respectable performance metrics, with precision, recall, F1 scores, and accuracy ranging from 62% to 67%. The Area Under the Curve (AUC) score of 68% indicated the model's effectiveness in making predictions.
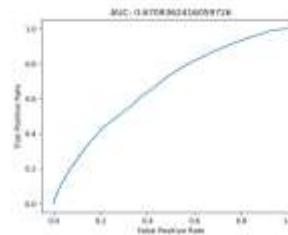


i.

III.  Hyper parameter Tuning and Validation

a.  Despite efforts to enhance the model through hyper parameter tuning using GridSearchCV, the accuracy score remained at 63%, consistent with the original model. Cross-validation and K-fold analysis validated this finding, yielding a consistent mean accuracy of 63%.

IV.     Tableau
        We created two tableau dashboards.



V.      Website Development

        a.  The website aims to provide valuable insights into homicide data from 2012 to 2017. It features an 'Arrest Prediction Page' offering predictive insights and a 'Dataset Overview Page' presenting a comprehensive view of the dataset. The website utilizes Python, Flask, Tableau, and was deployed on PythonAnywhere.  The website employs a dark color palette to maintain a tone of seriousness, respecting the gravity of the subject matter—homicide data representation. It prioritizes visual engagement while acknowledging the sensitive nature of the topic.

VI.     Conclusions

   o    The dataset exhibits a right-skewed distribution of victim ages, with a mean age of 31.8. The age range spans from newborns or babies (minimum age of 0) up to a maximum age of 102. Additionally, there are notable outliers in the data.
   o    The dataset reveals that the most prevalent age group among victims falls within the range of 19-29 years old. Following closely, the 30-44 age group is the next most common. Conversely, the 65+ age range exhibits the lowest frequency among the victim demographics.
   o    Although the LGBM model holds potential for predictions, achieving an accuracy of 63% may not meet the required level of reliability. It's crucial to acknowledge the dataset's limitations: it encompasses data from only 50 US states, spanning 2007-2017, with New York contributing only two years of data. With a more comprehensive and uniformly detailed dataset, we could anticipate more precise and dependable outcomes.