



# DATA-602 Final Project

Predicting Crime Types based on Historical Data:  
A Machine Learning Approach

By:  
Kalaifarasi Kaliappan

# Project Approach:

CRISP-DM( Cross-Industry Standard Process for Data Mining)



BU: Understand proj from business perspective.  
Understand the true goal

DU: Collect, describe and explore the data

DP: Select, clean, construct and integrate for the next step.

Modeling: Select, clean, create and Assess the model.

Evaluation: Evaluate and review the results and determine the next step.

Deployment: Plan deployment and monitoring.



## **Business Problem:**

The dataset Crime Data contains details regarding crimes that happened in various cities around the City of Los Angeles between January 2020 and the present day. The dataset has attributes for the location, date, and time of the offense as well as the type of crime that was committed. Based on the location, date, sex, race and time of the incident, a machine-learning model can be created using the dataset to forecast the sort of crime that will be committed.

The goal of this project is to create a machine-learning model that can identify the kind of crime being committed based on the location, the time and day it occurs, and other information. This can help law enforcement officials effectively utilize their resources and prevent crimes from happening in the future.



## About Dataset:

**Data Source:** <https://catalog.data.gov/dataset/crime-data-from-2020-to-present>

**Quantity:** 694,915 Records

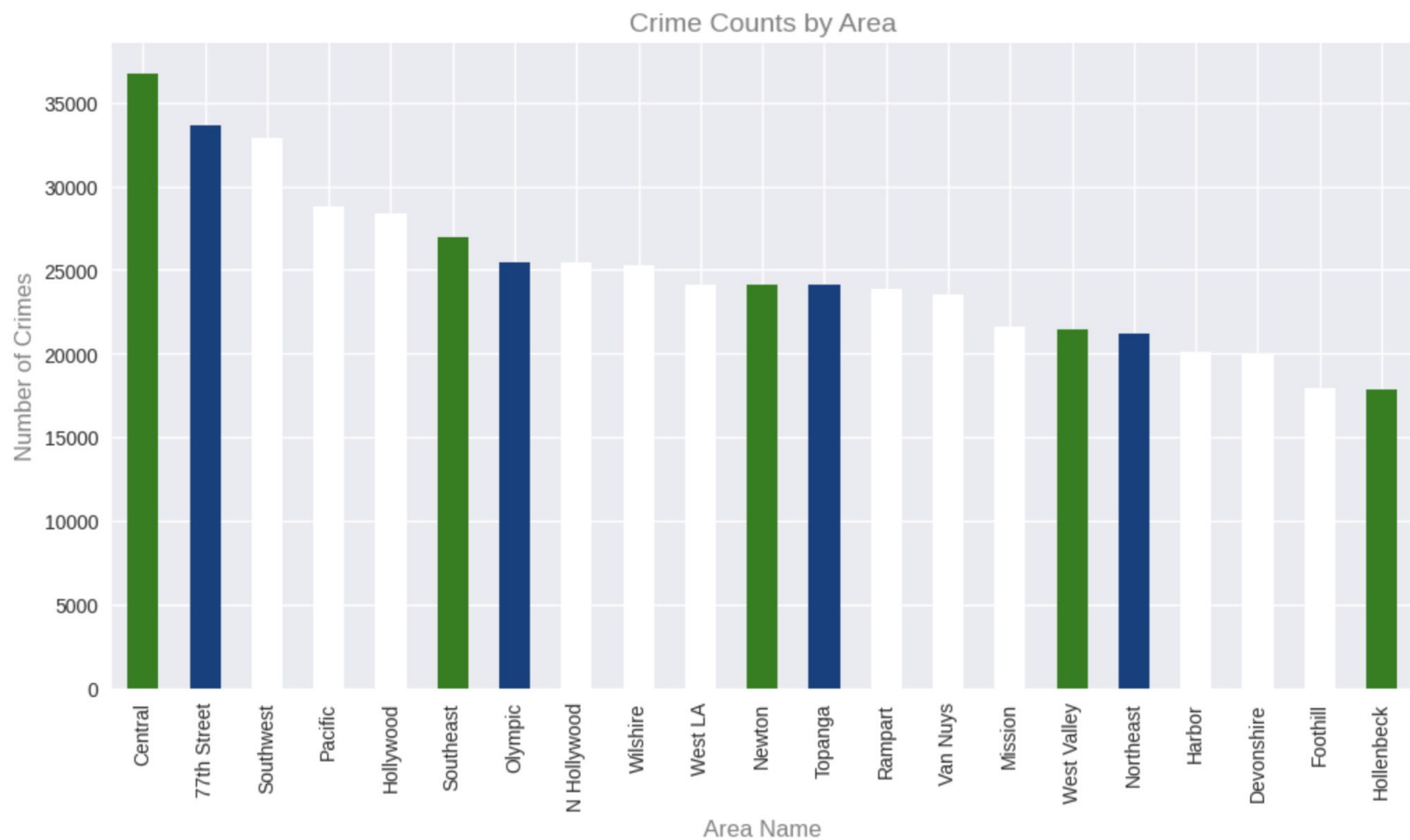
**Columns:** 28

**Key Attributes:** Location, Vict Age, Vict Sex, Vict Race, Data\_Occ, Crime description, Premis description.

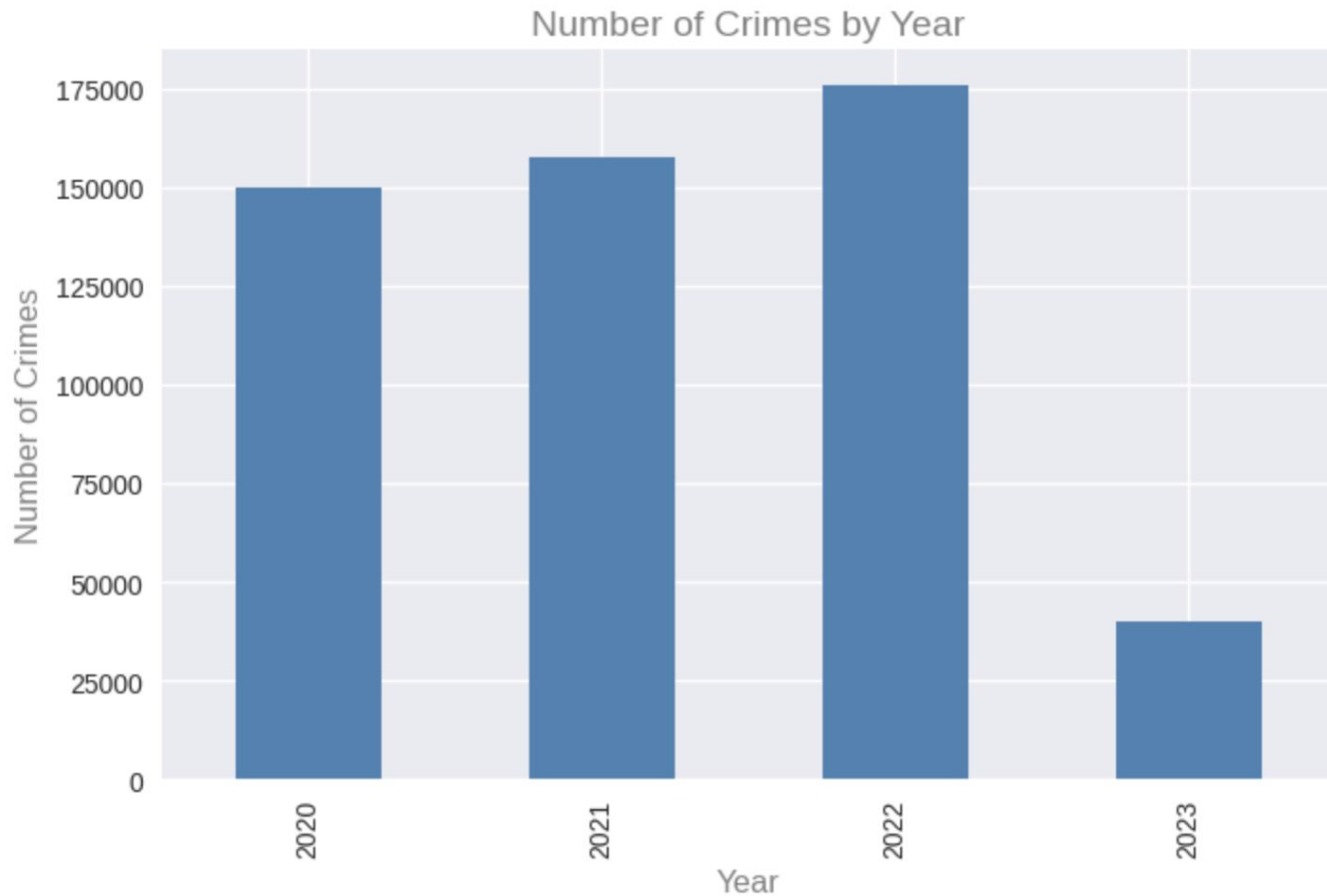


# Exploratory data analysis

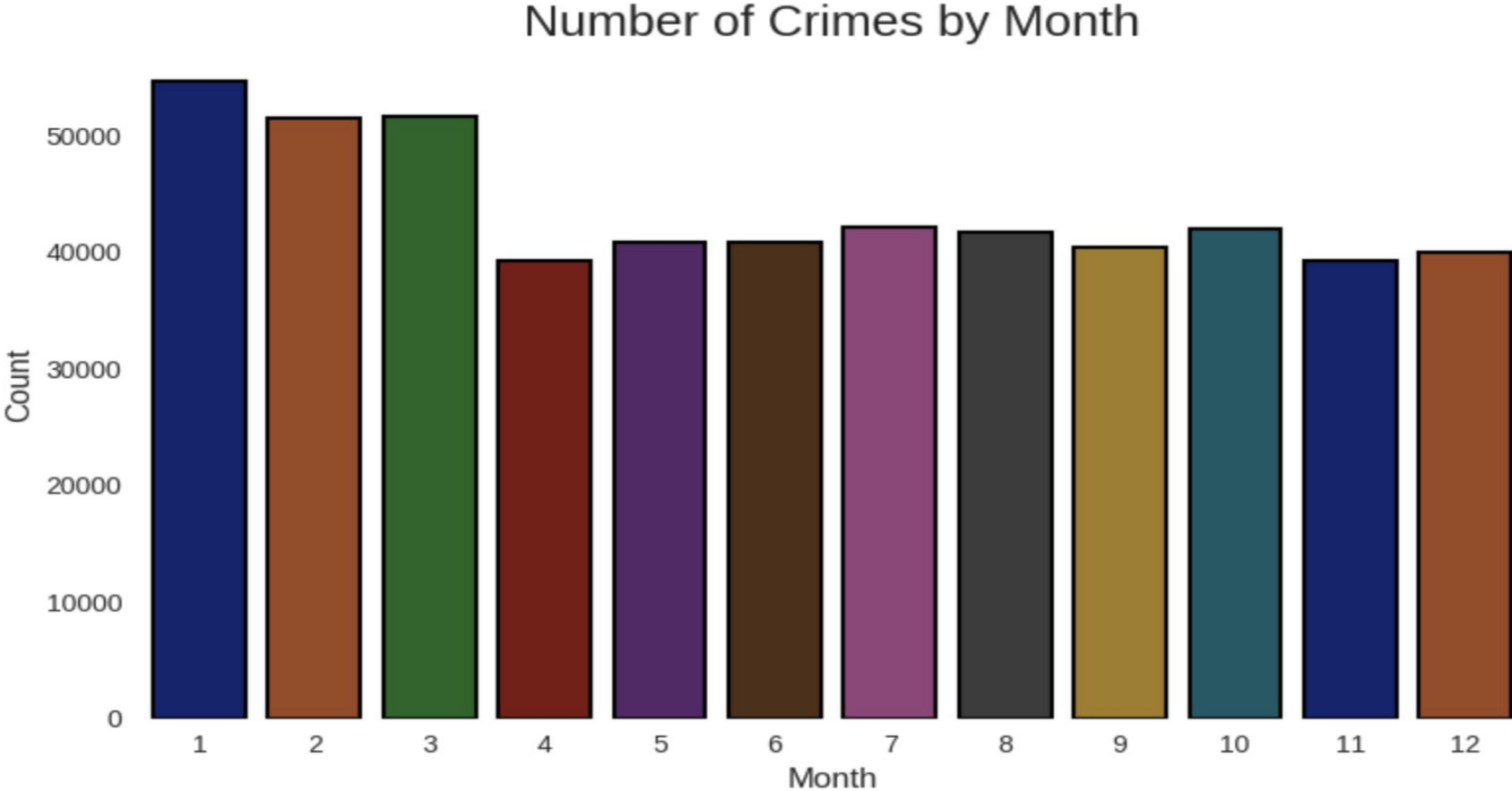
# Crime Counts By Area



# Number of Crimes by Year

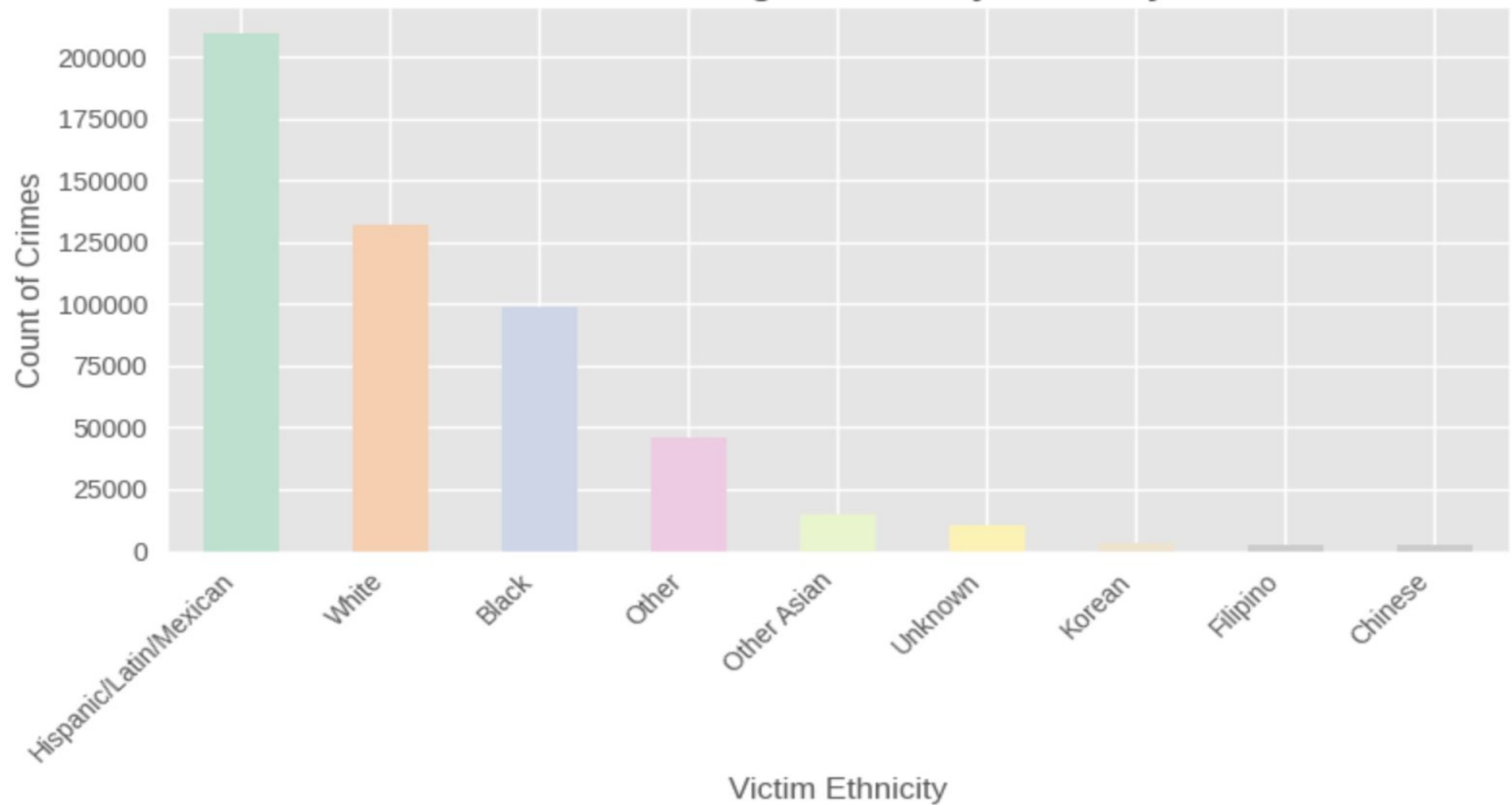


# Number of Crimes by Month:

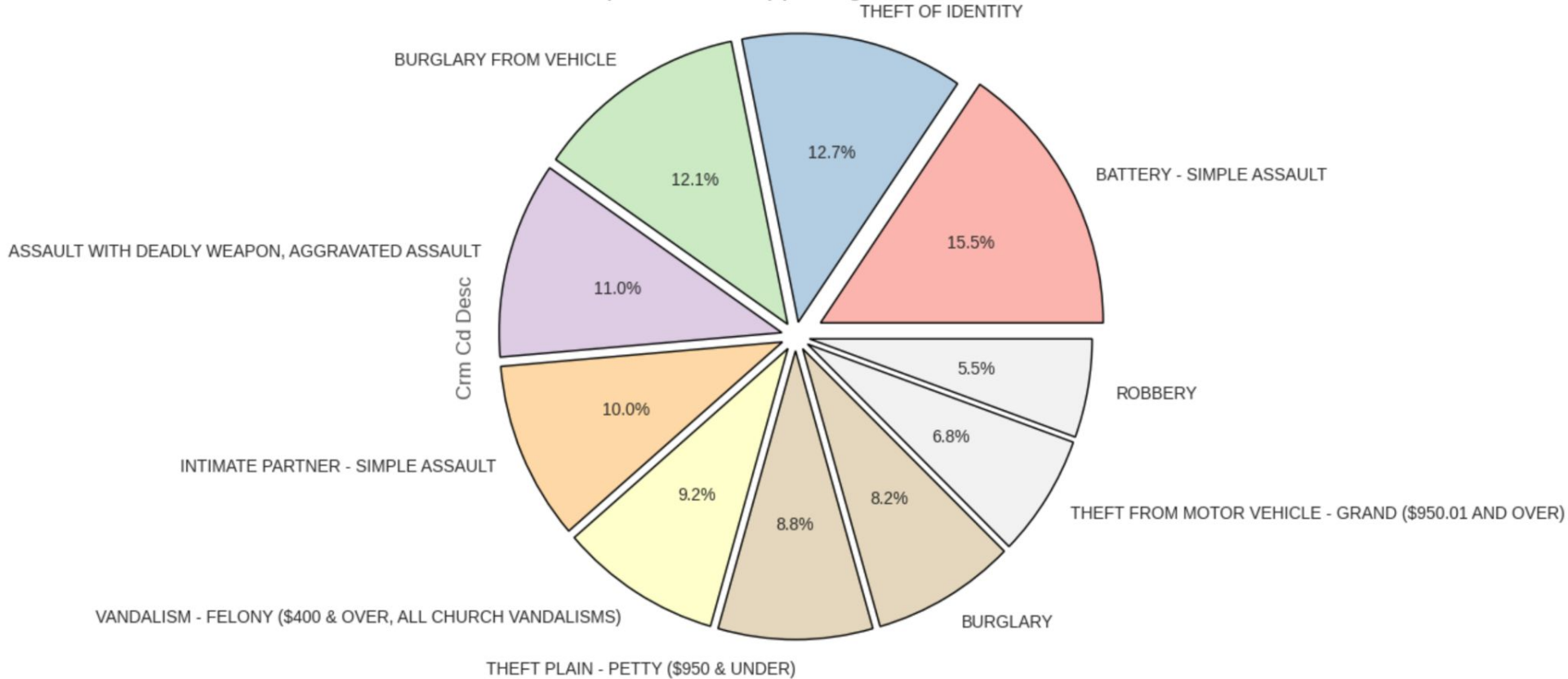




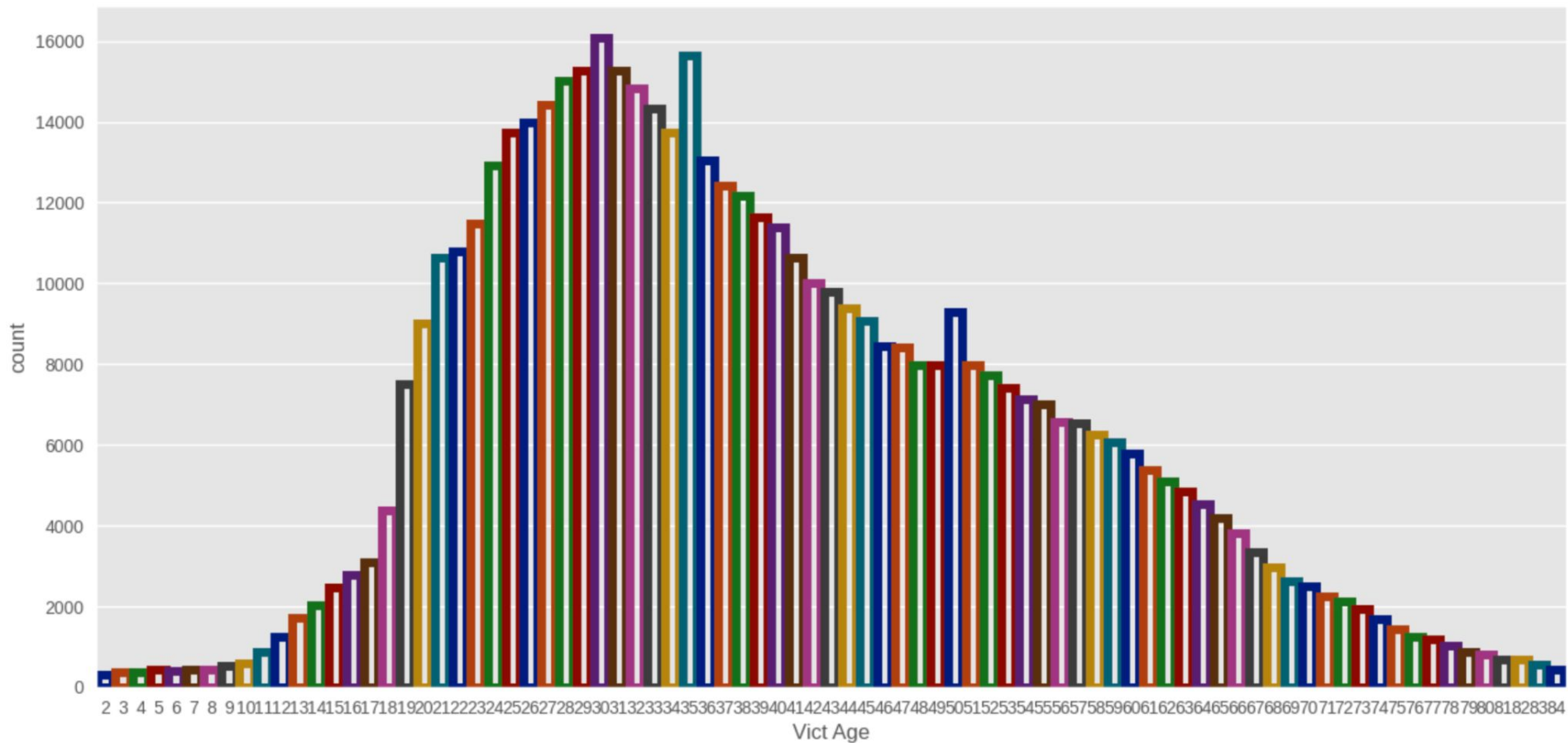
# Victims Categorization by Ethnicity



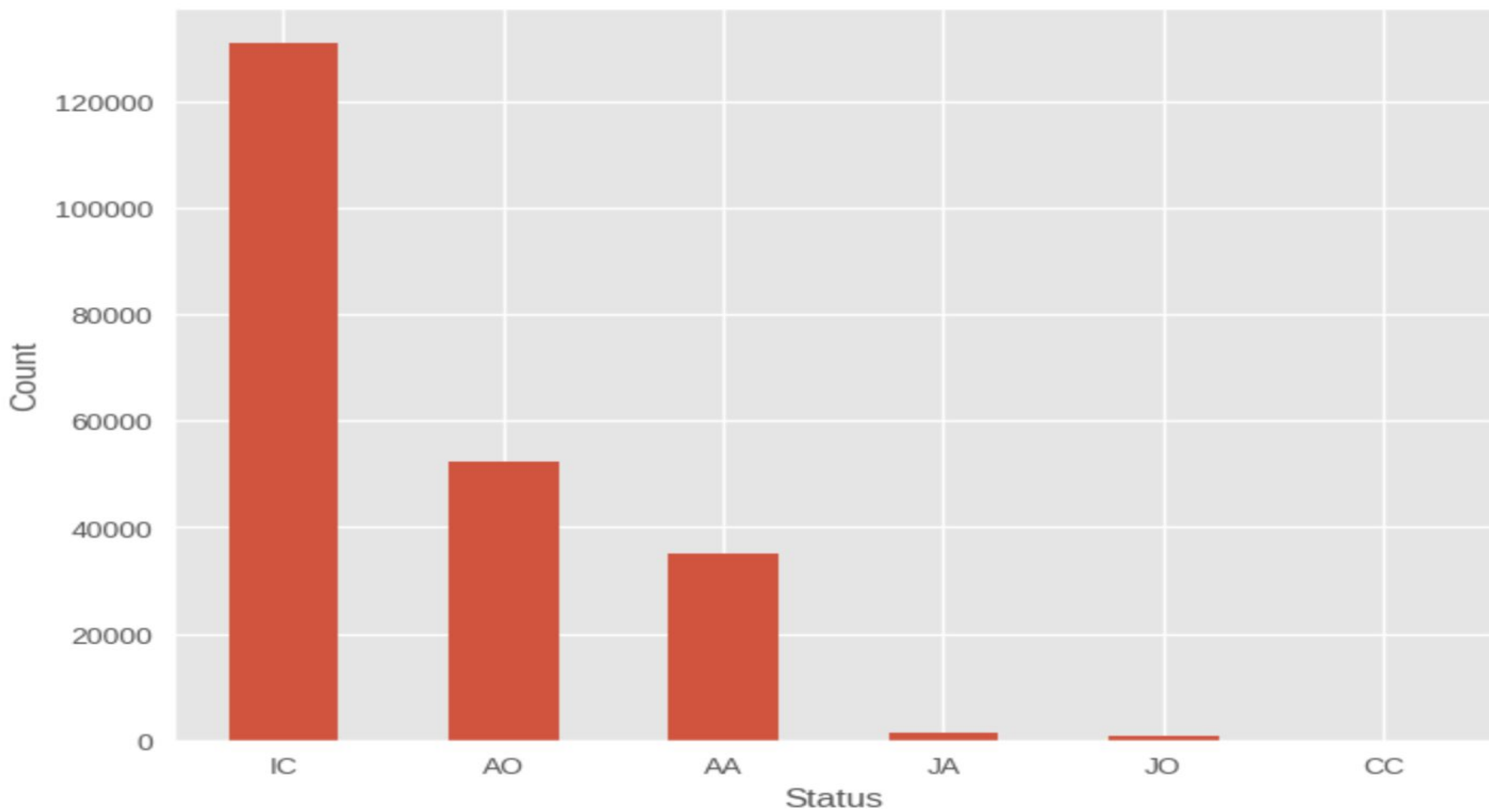
## Top 10 Most Happening Crimes in LA



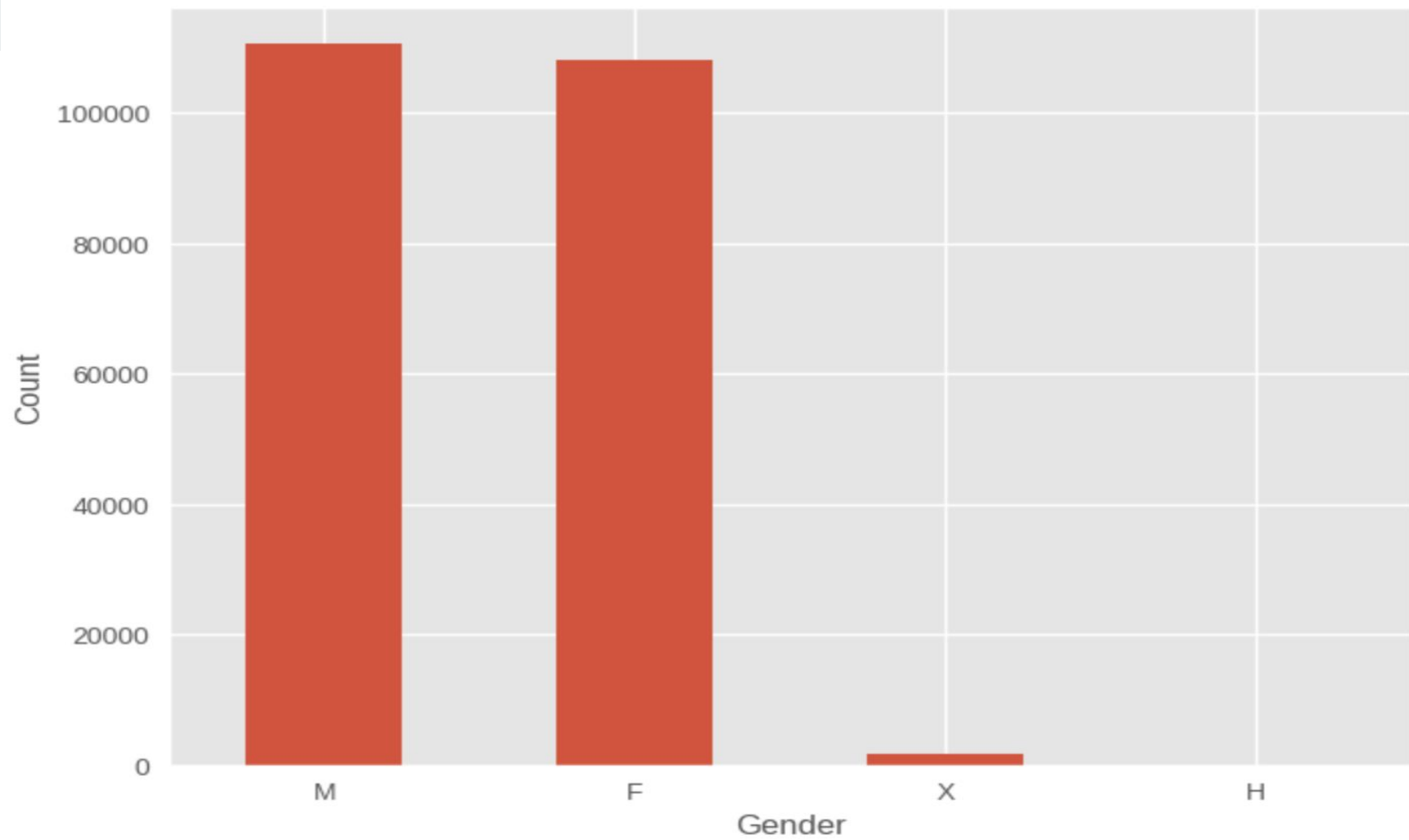
# Affected victim age distribution



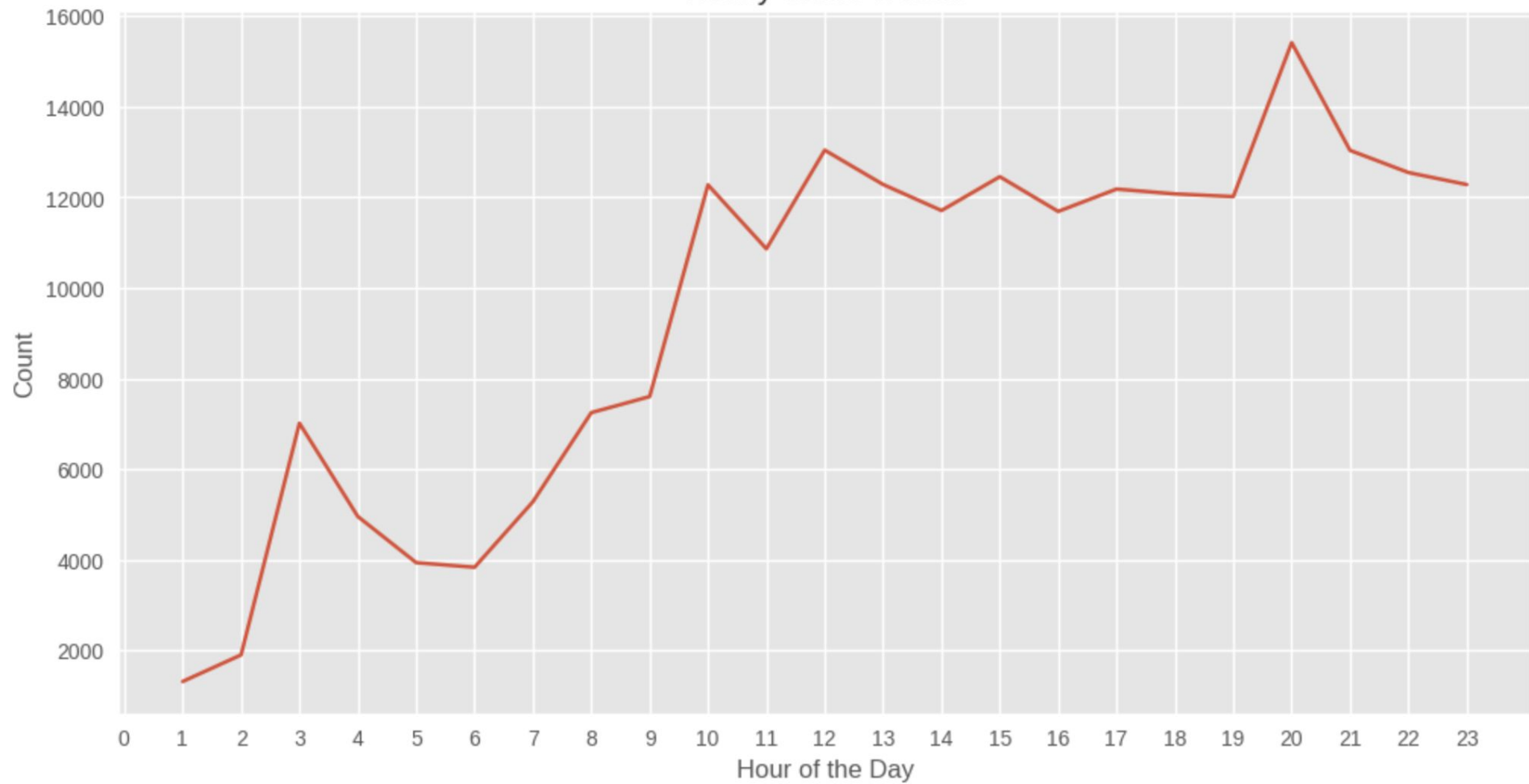
Distribution of Crime Status

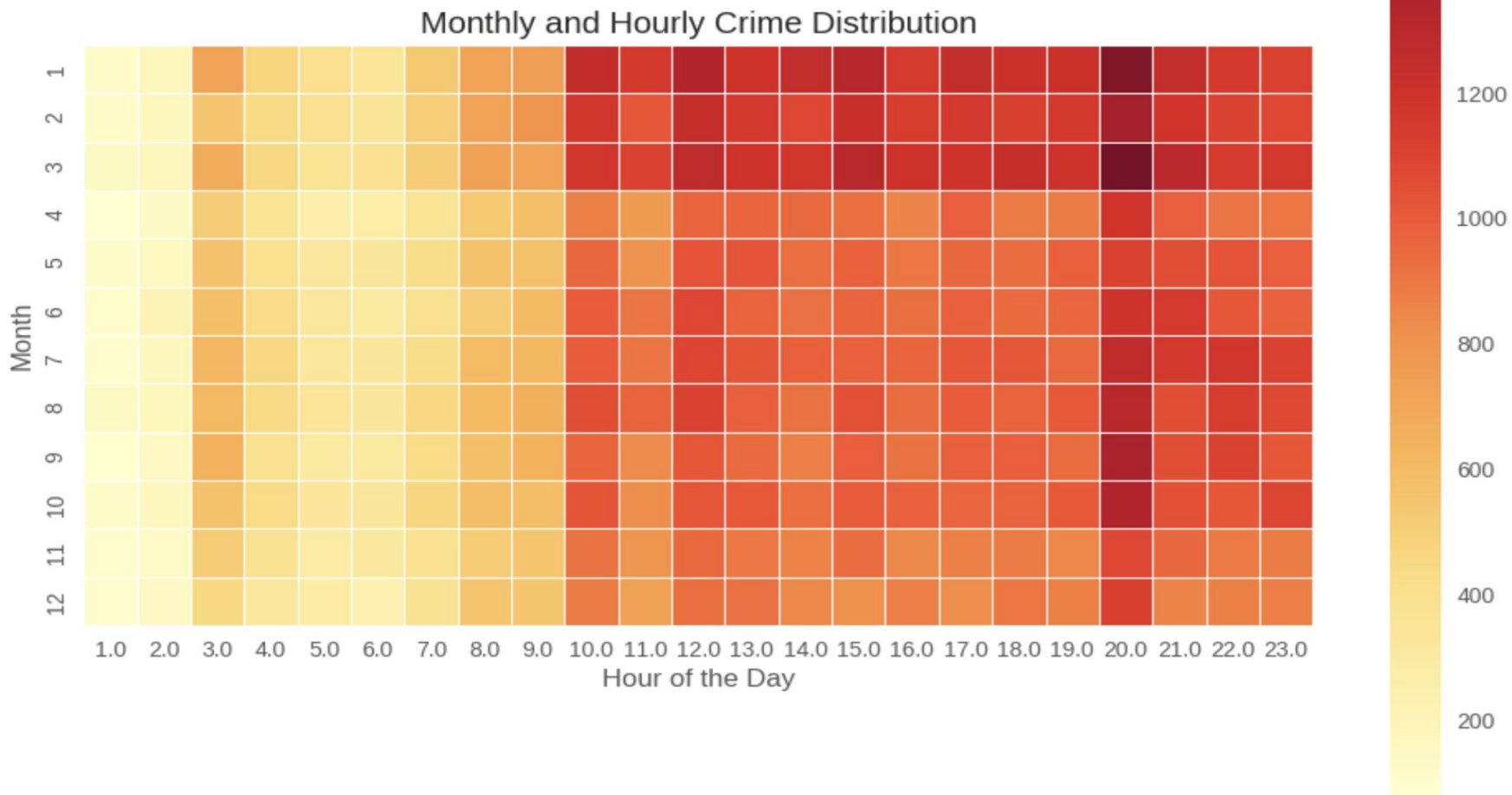


Distribution of Gender of Victims



# Hourly Crime Trends







# Modeling:

**Type of Modeling:** Classification

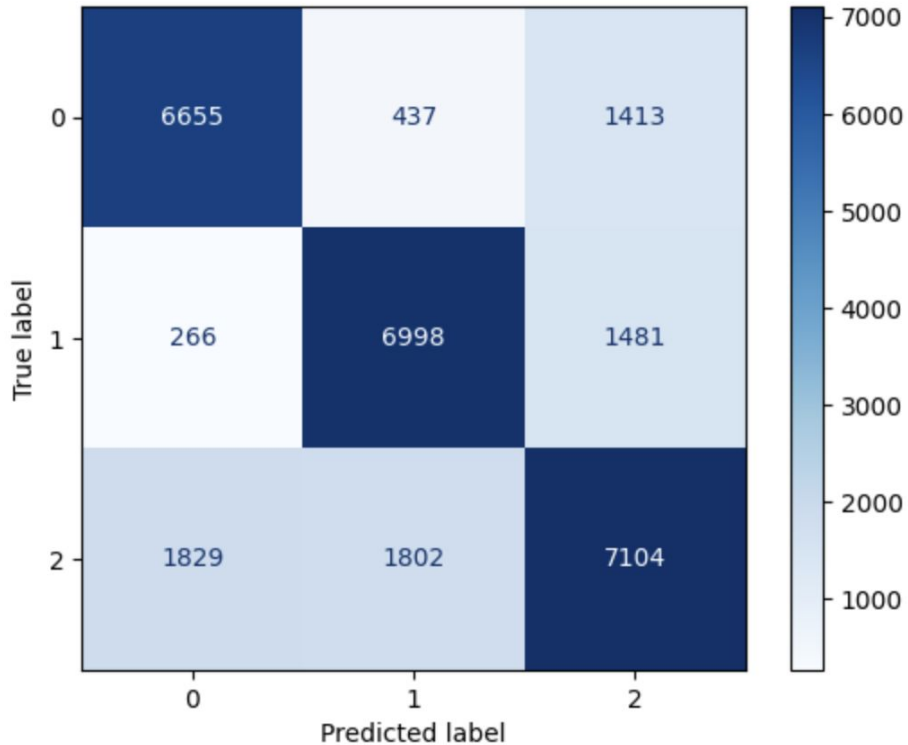
**Features:** Area, Vict\_Age, Vict\_Sex, Vict\_descent, Premis cd, Time\_occ,  
Day\_of\_week

**Target Variable:** crm cd



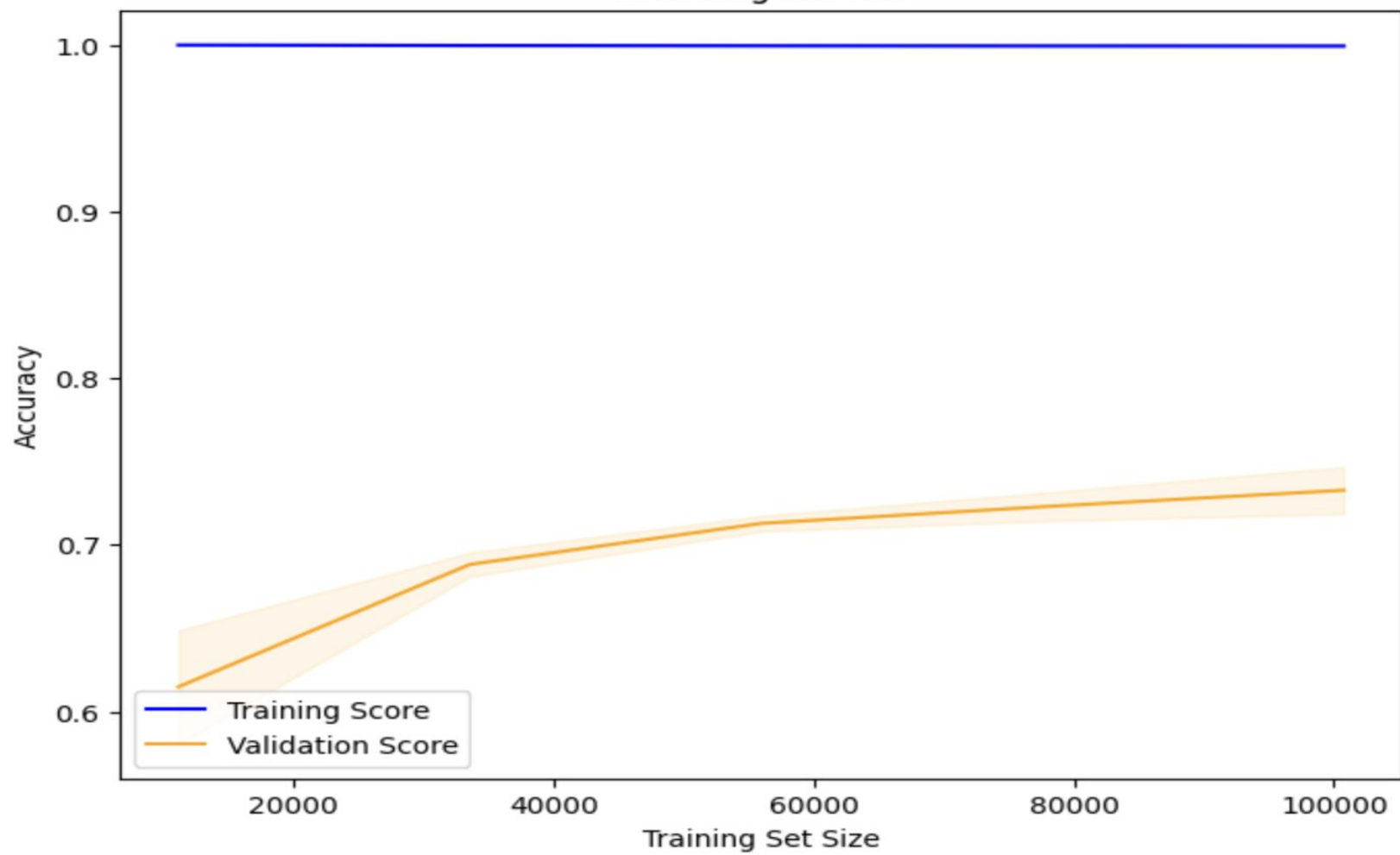
# Random Forest Classification:

Confusion Matrix

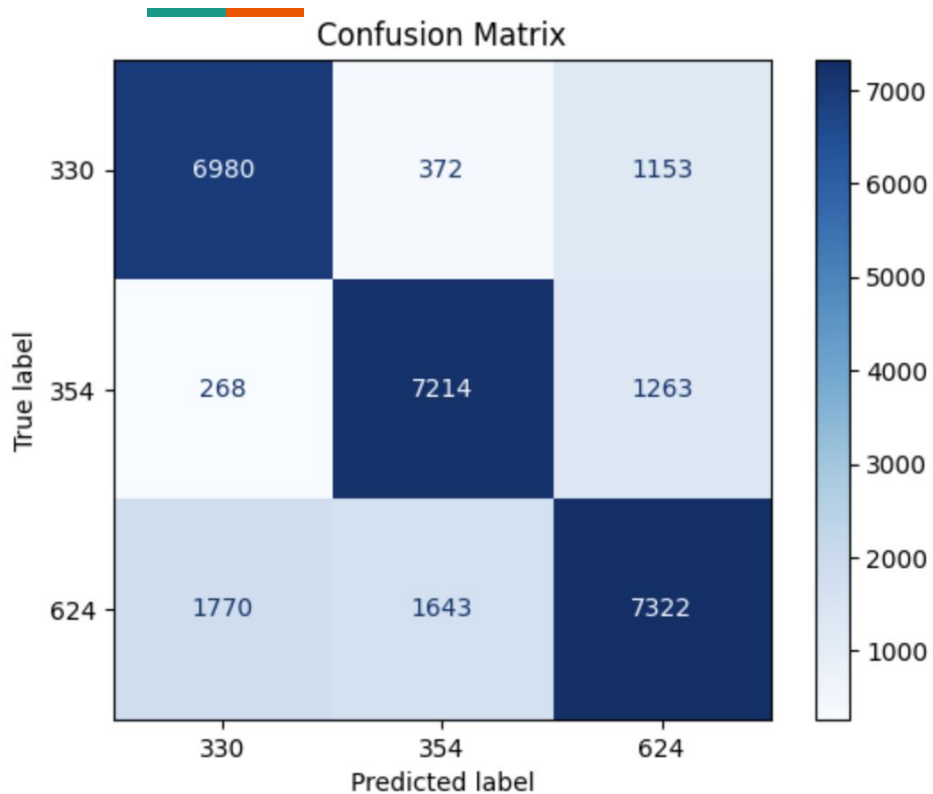


Accuracy: 0.742

Learning Curves

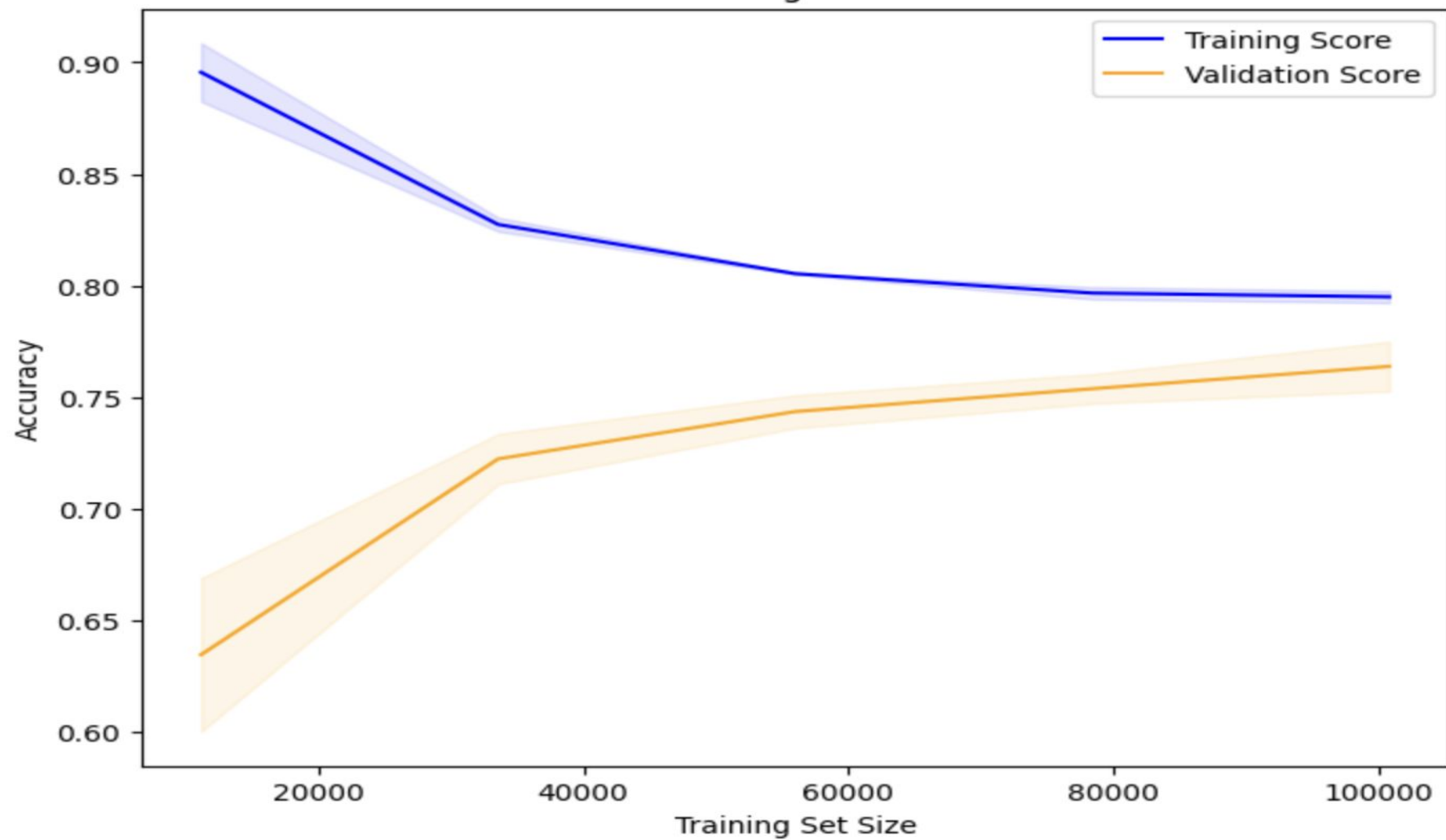


# XGBoost (Extreme Gradient Boosting)

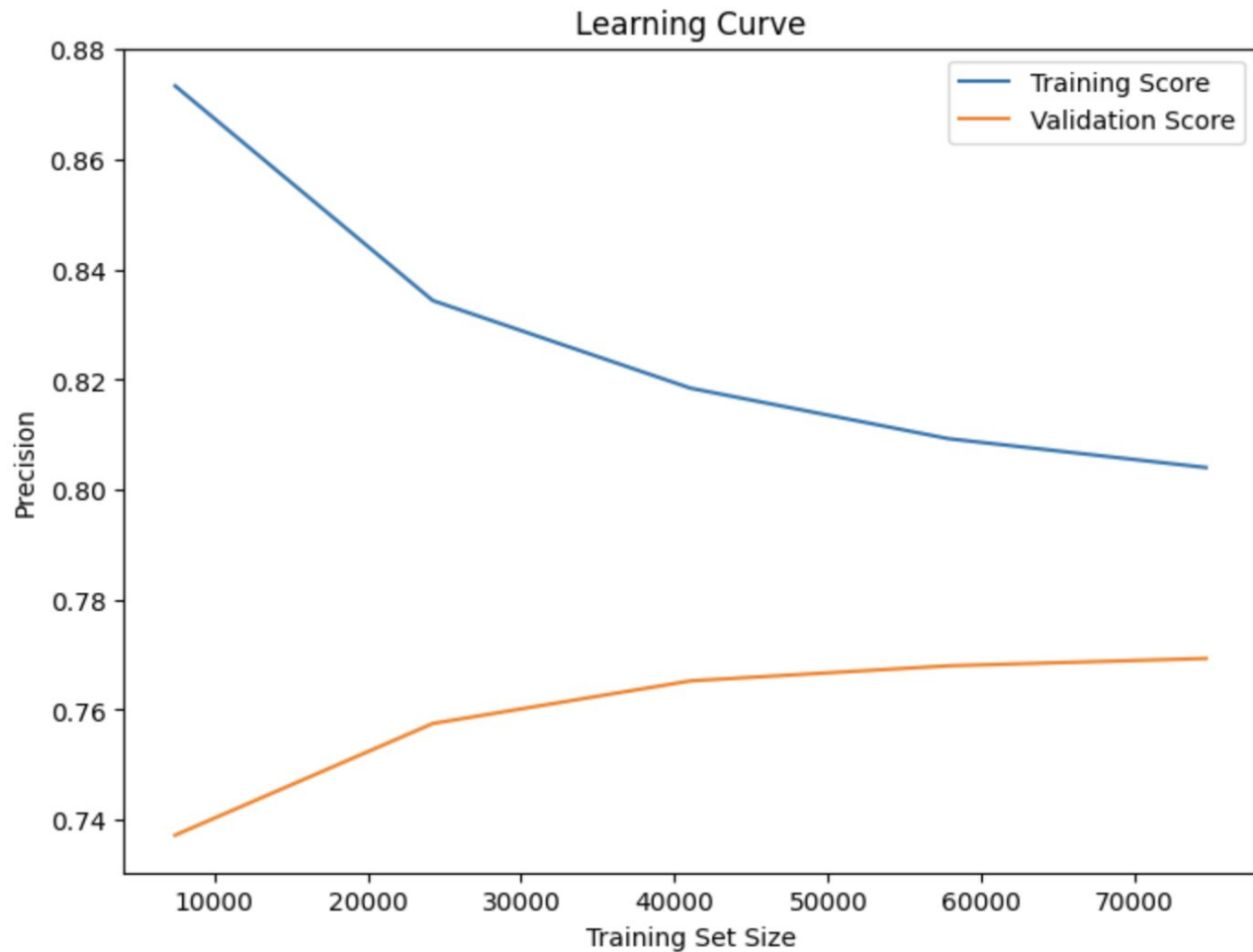


Accuracy: 0.76

# Learning Curves



# Stacking:



Precision:  
0.7701786972292215

# Advantages of Random Forest and XGboost model:



- **Ensemble Learning** (capture more complex relationships and reduce the impact of individual tree errors)
- **Handling Non-linearity** (handle non-linear relationships between features and the target variable)
- **Robust to Overfitting**
- **Feature Importance**
- **Handling Imbalanced Data**
- **Flexibility and Tunability**



**Thank You!!!**