# CRIME RATE PREDICTOR

**CLOUD COMPUTING**

AISHVWARYA IYER

STUDENT ID: 16227781

GAYATHREE IYER

STUDENT ID: 16227784

# 1. INTRODUCTION

## 1.1 MOTIVATION

According to The New York Times, America has definitely more violence when compared to other rich countries. Murder rates are also much higher in the United States than in countries like Japan, Europe, or Canada. We also have more assaults, and robberies than other rich countries. When compared to other affluent nations, crime rates have always been much higher in America. Hence predicting and analyzing crime data is extremely crucial.

Crime has also been the hot topic for the 'Smart City' conference.

## 1.2 GOAL

Crime Rate Predictor is a Cloud Computing project that focusses on applying machine learning algorithm to build models that will not only predict crime rates in various states but also predict the police funding and educational attainment required in the future to mitigate crime.
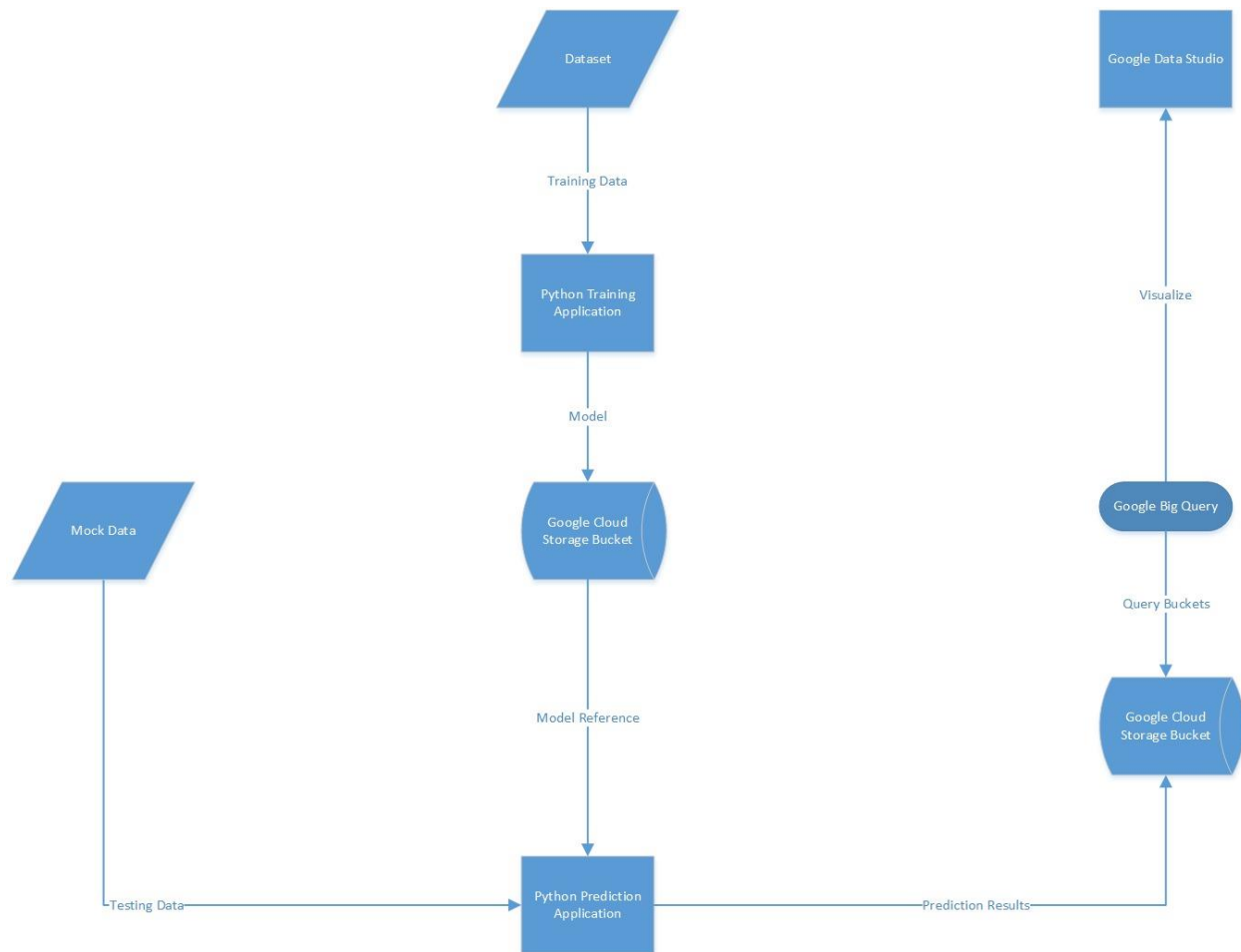
# 2. TEAM MEMBER'S CONTRIBUTION AREA

***Aishvwarya Iyer*:**

- Google Cloud Platform Infrastructure setup
- Python Training Program
- Google Big Query


***Gayathree Iyer:***

- Python Prediction Program
- PowerShell Automation Scripts
- Google Data Studio

## 3. SYSTEM ARCHITECTURE



**System Architecture**

## 4. MACHINE LEARNING ALGORITHM

We use **Multi variate linear regression model** for the prediction of crime. Multiple correlated dependent variables are predicted using this model. Regression analysis enables to describe the relationship between response variable and one or more predictor variables.

$$Yi = \alpha + \beta_1 x_i + \beta_2 x_i + \ldots + \beta_n x_i$$

$Y_i$ is the estimate of $i^{th}$ component of dependent variable y, where we have **n** independent variables and $x^j_i$ denotes the $i^{th}$ component of the $j^{th}$ independent variable/feature.

## 5. DATA SET DESCRIPTION

The data set contains a case study of education, crime, and police funding for small cities in ten southeastern and eastern states. The states are New York, South Carolina, Florida, Connecticut, Georgia, Rhode Island, New Hampshire, Maine, North Carolina and Virginia.


The data collected are for a sample of 50 small cities in these states.

X1 = Total Overall Reported Crime Rate per 1 Million Residents

X2 = Reported Violent Crime Rate per 100,000 Residents

X3 = Annual Police Funding (Dollars) per Resident

X4 = Percent of People 25 Years and Older That Have Had 4 years of High School

X5 = Percent of 16- to 19-Year-Olds Not in High School and Not High School Graduates

X6 = Percent of 18- to 24-Year-Olds Enrolled in College

X7 = Percent of People 25 Years and Older with at Least 4 Years of College

X8 = States


## 6. TECHNOLOGIES
- ➢ Programming language: Python.
- ➢ Google Cloud infrastructure is built using terraform scripts.
- ➢ Automation is done using PowerShell scripts.
- ➢ Source control management: Github

**Github link**: https://github.com/CloudComputing-Fall2107/CrimeRatePredictor

# 7. PROJECT MANAGEMENT

The work progress is tracked using Zenhub, which is an agile project management for Github.

## 8. PROJECT STEPS

```
CrimeRate.csv

(input)
```

PowerShell Script

Training  ──1──▶  training.py  ──2──▶  model1.pickle
                                        model2.pickle
                                        model3.pickle

Prediction  ──────▶  Read 3 Models
         ──4──

                        3

                Google
                Cloud Bucket

                     5

TestData1.csv
TestData2.csv  ──────▶  prediction.py
TestData3.csv

         6

         7

result1.csv

result2.csv

result3.csv

         8

Result
Google Cloud
Bucket

         9  ──────▶  Google Big Query

                        10

                Google Data
                Studio for
                visualization

## 8.1     INSTALLATION OF PACKAGES

Installation of necessary packages for the project using Choclatey, the package manager for windows.



**Install_Packages_PowerShell**

## 8.2     GOOGLE CLOUD PLATFORM INFRASTRUCTURE SETUP

Using Google Cloud PowerShell Tools.



**GCS Infratructure Setup**

## 8.3   TRAINING PROGRAM

The prediction strategy that we execute is to first train the features from training data and create training model using Pickle python library.

➢ **Predicting Total Overall Reported Crime:**

Here, we predict *the total overall reported crime rate per 1 million residents using:*

X2 = Reported Violent Crime Rate per 100,000 Residents

X3 = Annual Police Funding in Dollars per Resident

X4 = Percent of People 25 Years and Older That Have Had 4 years of high school

X5 = Percent of 16- to 19-Year-Olds Not in High School and not high school graduates

X6 = Percent of 18- to 24-Year-Olds Enrolled in college

X7 = Percent of People 25 Years and older with at Least 4 Years of college

X8 = States

➢ **Predicting Annual police Funding in Dollars:**

Here, we predict *the annual police funding in dollars per resident* using:

X1 = Total Overall Reported Crime Rate per 1 Million Residents

X2 = Reported Violent Crime Rate per 100,000 Residents

X8 = States

- Government can decide the annual police funding using this model efficiently to mitigate the crime rate.
- According to FBI statistics, annual police funding affects the crime rate to a great extent. Hence, it is very crucial to help the government know the annual police funding in advance.

➢ **Predicting the Age Group and Educational Attainments** :
Here, we predict the *age group and educational attainments* of the person committing the crime using:

X1 = Total Overall Reported Crime Rate per 1 Million Residents

X2 = Reported Violent Crime Rate per 100,000 Residents

X3 = Annual Police Funding in Dollars per Resident

X8 = States

- According to Alliance report findings, there is an indirect correlation between educational attainment and arrest rates.
- According to the most recent data from the U.S. Bureau of Justice, 56 percent of federal inmates, 67 percent of inmates in state prisons, and 69 percent of inmates in local jails did not complete high school.
- Plus, when examining total crime savings, the report also forecasts the number of individual crimes that could be avoided by expanding the high school graduation rate by 5 percent, and finds that such an increase would decrease overall annual incidences of larceny by more than 37k; assault by nearly 60k; burglaries by more than 17k and motor vehicle theft by more than 31k. It would also avert nearly 1,300 murders, more than 1,500 robberies and more than 3,800 occurrences of rape.



**Training_Python**

## 8.4    PREDICTION PROGRAM

Using the model we predict necessary features.



**Prediction_Python**

## 8.5    PROJECT AUTOMATION SCRIPTS

Automation is done using PowerShell scripts.
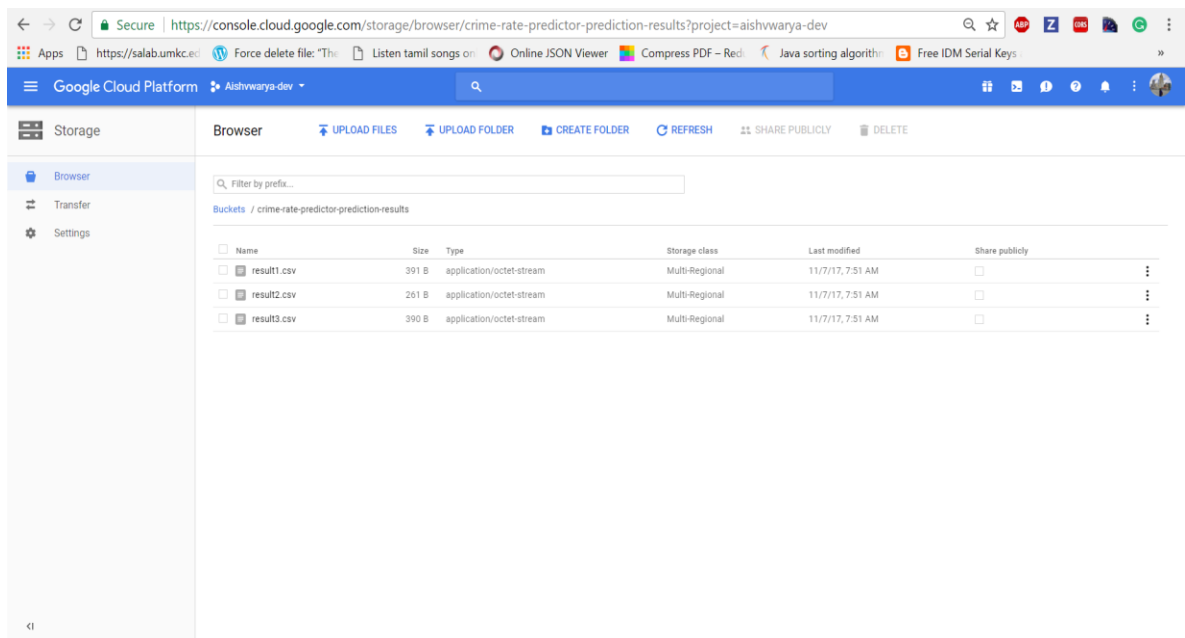


**Automation_PowerShell**

## 8.6    GOOGLE STORAGE BUCKET

We store the necessary model and the prediction results in appropriate buckets.


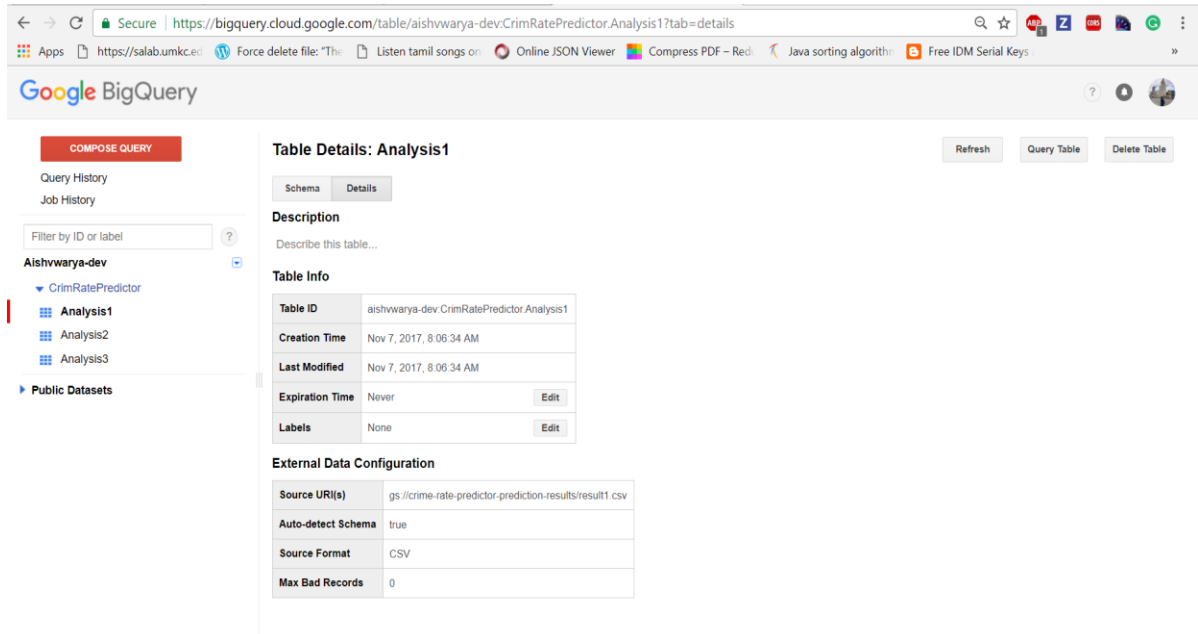
**GCS- Bucket Modules**



**GCS- Bucket Results**

## 8.7   GOOGLE BIG QUERY

We created Google Big Query tables from source which is google cloud storage bucket files. After the table is created, we can see schema of csv files and its configurations. We use Google Big Query to query the CSV files stored in GCS buckets.



**GCS Big Query Details**



**GCS Big Query Schema**

## 8.8    GOOGLE DATA STUDIO

Visualize big query results in data studio by choosing data source as BigQuery table.

- We use Google data studio which is a business analytics tool to visualize the prediction results.
- First, we create data source from big query and choose appropriate table.
- After creating data source, we build metrics like Overall expected crimes by calculating sum of all expected crimes.
- Then, we create pie chart, geo locations and bar chart for visualizing prediction results.



**Data Studio HomePage**

**Data Studio Choosing Data Source**



**Data Studio Data Source Schema**

**Data Studio Charts Sum**



**Data Studio Charts Edit**

## 8.9    RESULTS

The geo-location chart result:

- Here, we visualize the expected crime per region. More the intensity of color red, more the crime.
- The result shows '**Georgia**' with expected crime of 880 and overall expected crime of all regions is 6,900 approximately.

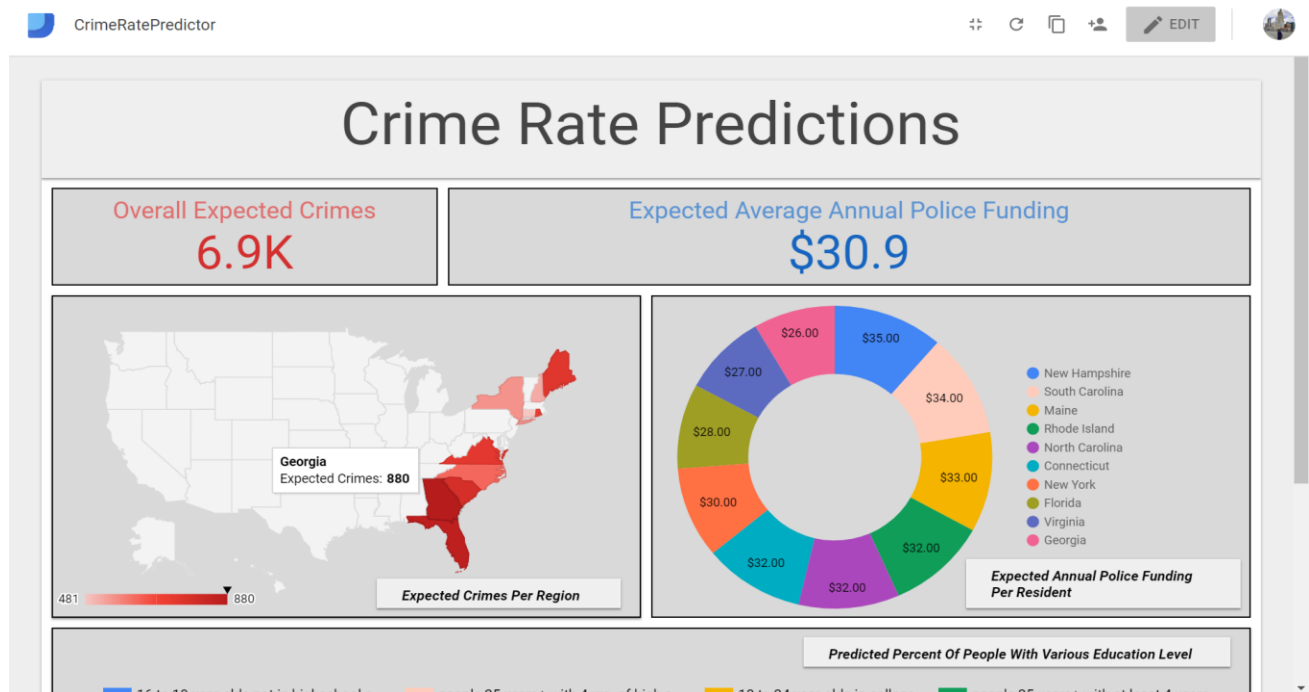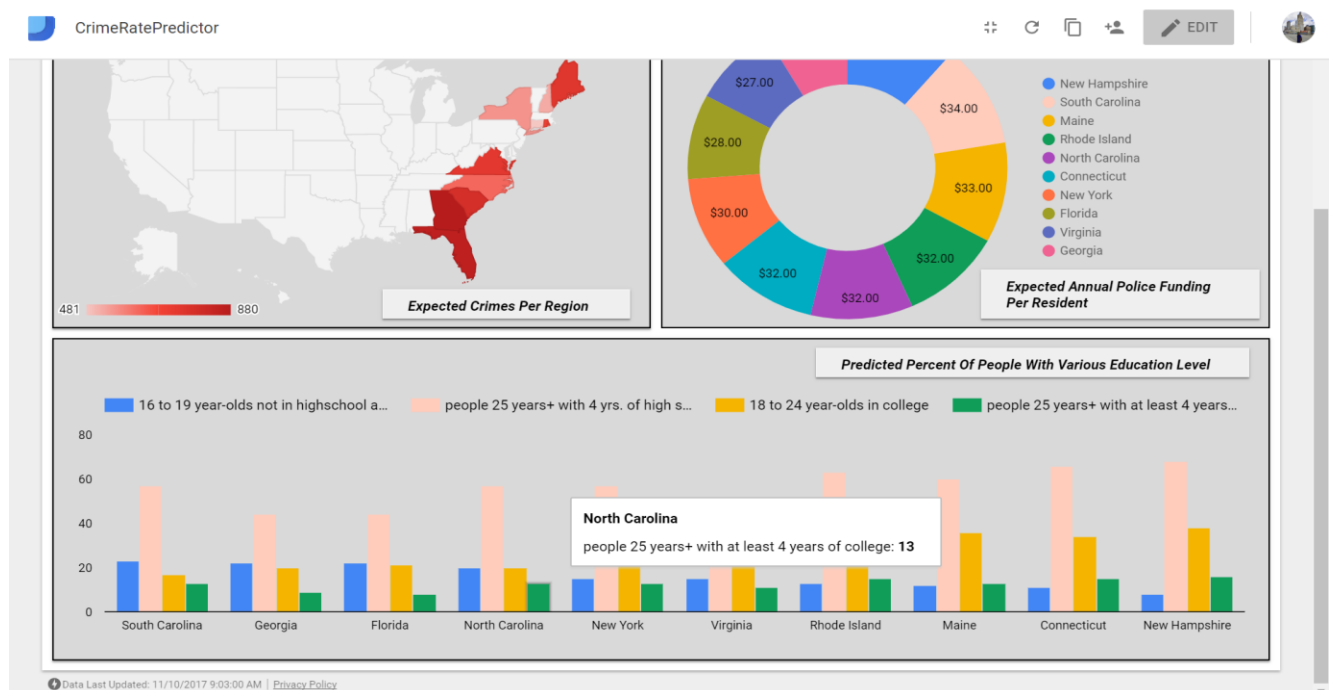The pie chart result:

- Here, we visualize expected average annual police funding.
- The region '**New Hampshire**' requires most of the police funding with 11.3% of total funding.

The bar chart result:

- Here, we visualize the predicted percentage of people with various educational level.
- We see that **Georgia** and **Florida** have less literacy rate compared to other regions.
- This also reflects in **Expected Crime per region** where the intensity of red (crime) is maximum.
- Thus, our prediction algorithm works with maximum accuracy.

**Data Studio Result-1**



**Data Studio Result-2**

**Link** - https://datastudio.google.com/open/1oq4-lb5RMd2S1crfN9MFS7gFHBoTMJEr

## 9. REFERENCE

- http://college.cengage.com/mathematics/brase/understandable_statistics/7e/students/datasets/mlr/frames/frame.html
- https://www.brookings.edu/research/more-cops/
- https://all4ed.org/press/crime-rates-linked-to-educational-attainment-new-alliance-report-finds/
- https://www.hackerearth.com/practice/machine-learning/linear-regression/multivariate-linear-regression-1/tutorial/
- http://ptl.sys.virginia.edu/ptl/sites/default/files/Area-Specific%20Crime%20Prediction%20Models.pdf