

# Insights for Chicago to Improve Public Safety

## Implementation Report

Bekalu Tadesse | CSCI622

### 1. Introduction

This project aims to support data-driven public safety initiatives for Chicago by applying data engineering principles, leveraging public datasets from data.gov, cloud and technology.

Through data ingestion, transformation, and analysis, the project aims to uncover valuable insights that can contribute to improving public safety.

The analysis seeks to provide actionable insights to policymakers, enabling them to identify specific geographic areas that require attention and make informed decisions to enhance public safety effectively.

### 2. Project Setup and Architecture

#### 2.1 Overall architecture

The project ingests 4 datasets (2 dynamic and 2 static) from data.gov API's into azure cloud blob storage (object storage). Ingestion is carried out using a python code that runs in a pip env virtual environment on a local machine.

Exploratory data analysis and Transformation is done within Azure databricks to unify datasets, merge and aggregate data. The ingestion code as well as Azure databricks transformation code is synced to a GitHub repository.

The final step of analysis/serving is done using PowerBI and Jupyter Notebook.

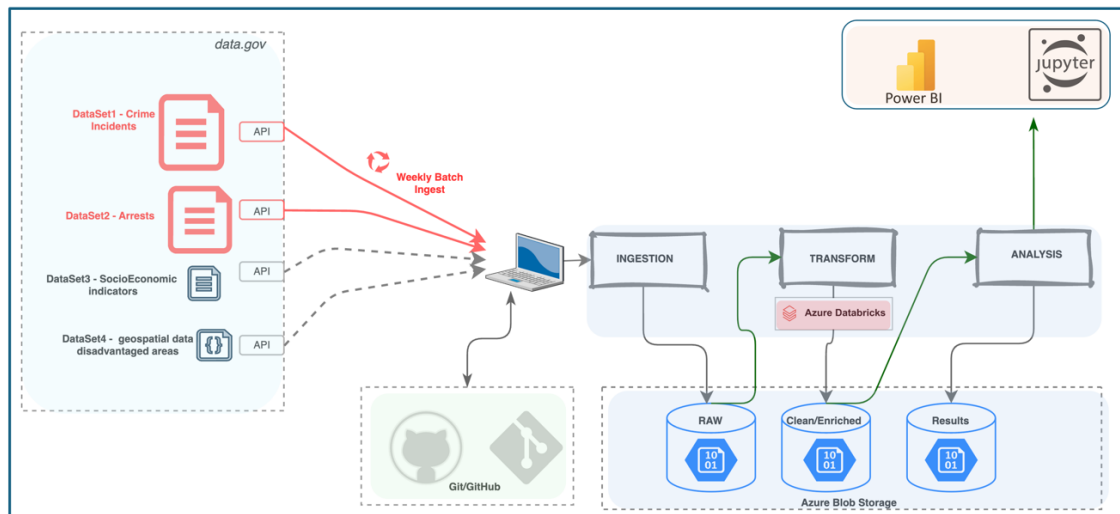


Figure 1 - Overall architecture

## Azure Cloud Setup

Within azure cloud a storage account, key vault and Azure data bricks workspace is created and configured using point and click. Automation of this task will be done in the future to improve this project further.

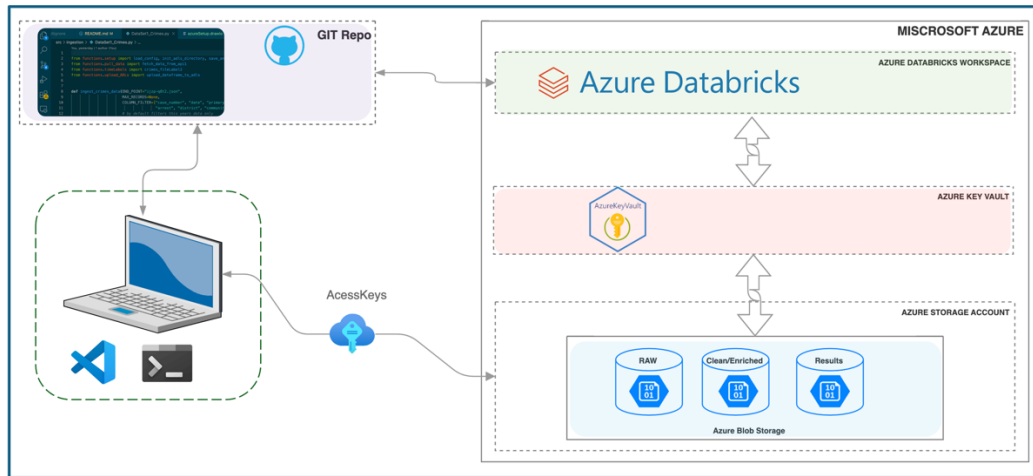


Figure 2 - Azure Cloud Setup

Three Blob Containers—designated for raw data, transformed data, and results—are **programmatically** created within the Azure storage account using Python code responsible for data ingestion. During the ingestion process, data is first retrieved from the API, saved as CSV files on a local machine as an intermediate step, and then uploaded to the respective containers in the Azure storage account.

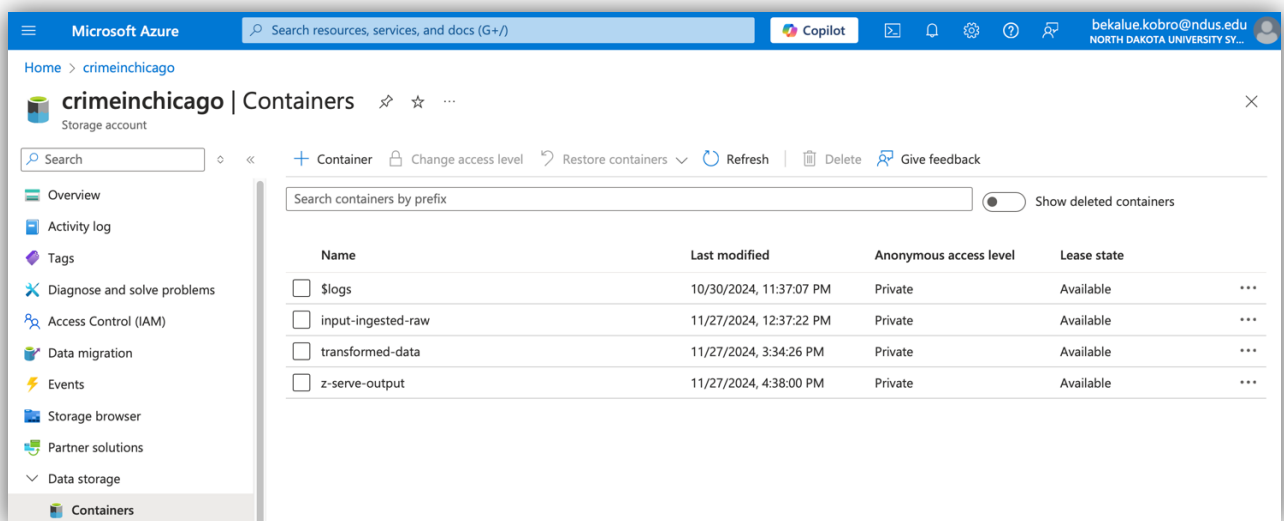


Figure 3 - blob containers

### Local machine setup (virtual environment)

A **Pipenv** virtual environment has been created on the local machine to run the ingestion code. The necessary Python packages are installed within this environment to ensure all dependencies are met for the process.

The **Pipfile** or **Pipfile.lock** can be used to recreate the development environment. Some of python packages installed and their use case is shown on the table below.

Python package	Use Case
azure-storage-file-datalake	- Programmatic ingestion to azure blob storage
tqdm	- Show progress bar during ingestion
Folium, geopandas	- visualize geospatial data
Plotly, matplotlib	- visualizations within jupyter notebook

## 3. Implementation Steps (Road map)

The implementation follows the data engineering lifecycle, consisting of three main steps: ingestion, transformation, and analysis/serving.

During the **ingestion** phase, data is retrieved from the Data.gov API and stored in Azure Blob Storage. In the **transformation** phase, the data is unified, merged, and aggregated to prepare it for the next step. In the **analysis** phase, the transformed data is used for analysis in Power BI and Jupyter Notebook.

In **Jupyter Notebook**, the aggregated and merged data is used to create visualizations, while in **Power BI**, the merged data is utilized to build dashboards.

### 3.1 Ingestion

#### 3.1.1 Data sources:

Four data sources has been used in this project. These are:

DataSource1: Crimes	<ul style="list-style-type: none"><li>- Dynamic data source</li><li>- updated frequently</li><li>- shows crime incidents</li><li>- 8.19 rows x 22 columns</li></ul>
DataSource2: Arrests	<ul style="list-style-type: none"><li>- Dynamic data source</li></ul>

	<ul style="list-style-type: none"> <li>- updated frequently</li> <li>- shows arrests</li> <li>- 660k rows x 24 rows</li> </ul>
DataSource3: Socio-economic indicators	<ul style="list-style-type: none"> <li>- Static data source</li> <li>- updated less frequently</li> <li>- shows economic indicators such as education, employment, poverty, etc.</li> <li>- 78 rows x 9 columns</li> </ul>
DataSource4: Socioeconomically disadvantaged areas	<ul style="list-style-type: none"> <li>- Static data source (geospatial data)</li> <li>- updated less frequently</li> <li>- shows areas in Chicago that have poor socio economic indicators</li> <li>- 5254 rows x 1 column</li> </ul>

### 3.2 Ingestion initial steps:

#### Create API Key & Azure Access:

- Register on Data.gov and create an API key using the "SignUpforAppToken" option.
- Store the API key in a config file.
- Also copy SAS key into another config file to be able to write to azure storage programmatically.

#### Data Retrieval via API

- Pull data via the Data.gov API.
- **Libraries Required:** Install pandas, soapy, jupyter, and azure-storage-file-datalake.

#### Programmatic Storage

- Create a directory within a container in Azure Blob Storage programmatically.
- Store datasets in .csv format for easy access in the next step.

#### 3.2.1 Dynamic data Ingestion: Batch and row/column Filter

A batch ingestion technique is employed to handle dynamic data sources. Since these sources include a significant number of records (over 8 million, and 600k in data sets 1 and 2 respectively) a row and column filters are applied at the API level. This approach prevents API throttling and reduces the need for extensive data cleaning later.

Column filters are used to retrieve only the relevant columns, while additional row filters are incorporated into the ingestion code to exclude null records and process the data in smaller, manageable chunks, organized by year as shown on the figures below.

Microsoft Azure | Search resources, services, and docs (G+ /)

Home > crimeinchicago | Containers >

**input-ingested-raw** Container

Authentication method: Access key (Switch to Microsoft Entra user account)  
Location: input-ingested-raw / Arrests

Search blobs by prefix (case-sensitive)  ☐ Show deleted objects

Name	Modified	Access tier	Archive status	Blob type	Size
[.]					
Arrests_2019-01-01_to_2019-12-31_50753_rows.csv	11/27/2024, 3:02:10 P...	Hot (Inferred)		Block blob	3.59 I
Arrests_2020-01-01_to_2020-12-31_31354_rows.csv	11/27/2024, 3:04:41 P...	Hot (Inferred)		Block blob	2.22 I
Arrests_2021-01-01_to_2021-12-31_24422_rows.csv	11/27/2024, 3:06:59 P...	Hot (Inferred)		Block blob	1.72 I
Arrests_2022-01-01_to_2022-12-31_25177_rows.csv	11/27/2024, 3:09:02 P...	Hot (Inferred)		Block blob	1.76 I
Arrests_2023-01-01_to_2023-12-31_28990_rows.csv	11/27/2024, 3:11:25 P...	Hot (Inferred)		Block blob	2.03 I
Arrests_2024-01-01_to_2024-11-23_29945_rows.csv	11/27/2024, 3:13:49 P...	Hot (Inferred)		Block blob	2.12 I

Microsoft Azure | Search resources, services, and docs (G+ /)

Home > crimeinchicago | Containers >

**input-ingested-raw** Container

Authentication method: Access key (Switch to Microsoft Entra user account)  
Location: input-ingested-raw / Crime2019\_to\_Present

Search blobs by prefix (case-sensitive)  ☐ Show deleted objects

Name	Modified	Access tier	Archive status	Blob type	Size
[.]					
Crimes_2019-01-01_to_2019-12-31_56106_rows.csv	11/27/2024, 2:34:02 P...	Hot (Inferred)		Block blob	5.96 I
Crimes_2020-01-01_to_2020-12-31_33808_rows.csv	11/27/2024, 2:36:44 P...	Hot (Inferred)		Block blob	3.67 I
Crimes_2021-01-01_to_2021-12-31_25095_rows.csv	11/27/2024, 2:38:48 P...	Hot (Inferred)		Block blob	2.79 I
Crimes_2022-01-01_to_2022-12-31_27248_rows.csv	11/27/2024, 2:41:02 P...	Hot (Inferred)		Block blob	3.02 I
Crimes_2023-01-01_to_2023-12-31_31560_rows.csv	11/27/2024, 2:43:38 P...	Hot (Inferred)		Block blob	3.49 I
Crimes_2024-01-01_to_2024-11-19_30397_rows.csv	11/27/2024, 2:46:06 P...	Hot (Inferred)		Block blob	3.34 I

Figure 4 - chunks of data by year after batch ingestion by year

After the historical data is fully retrieved, continuous ingestion is performed by applying row filters with a specified time window. This allows for the dynamic datasets to be ingested weekly as new data becomes available.

### 3.3 Transformation: Unify, Join, Aggregate data

Data transformation involves unifying the segmented data from each source into unified datasets. These unified datasets are then merged into a single large table using common attributes/columns as the basis for the joins.

Summary of joining attributes for the data sets is listed on the table below:

DataSets	Join_attribute
DataSet1_Crimes & DataSet2_Arrests	case_number
Dataset_2_Arrests & DataSet3_incicators	Community_area_number
Dataset1_crimes & Dataset4_disadvantagedAreas	Latitude, longitude

### Dataset Aggregation:

The merged data is further aggregated based on specific indicator columns to enable deeper analysis within Jupyter Notebook.

### Simple Data Transformations with PowerBI:

Simpler data transformations are also carried out within power BI such as adding a new column, and some aggregations using DAX.

### Data Lineage:

After ingestion of Historical dynamic data in yearly batches, these yearly segments of data are then consolidated during the transformation phase to produce a unified CSV file. Finally, datasets are joined by common attributes as shown on the high-level data lineage graph below.

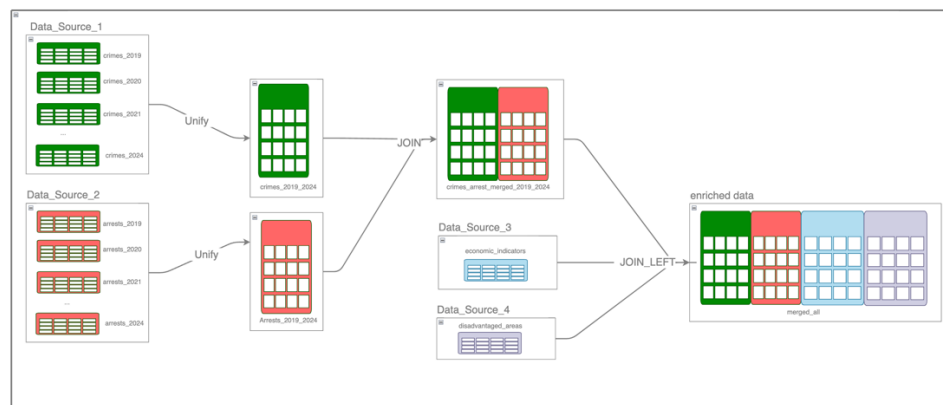


Figure 5 - Data lineage shows changes made to the data after ingestion down to transformation steps

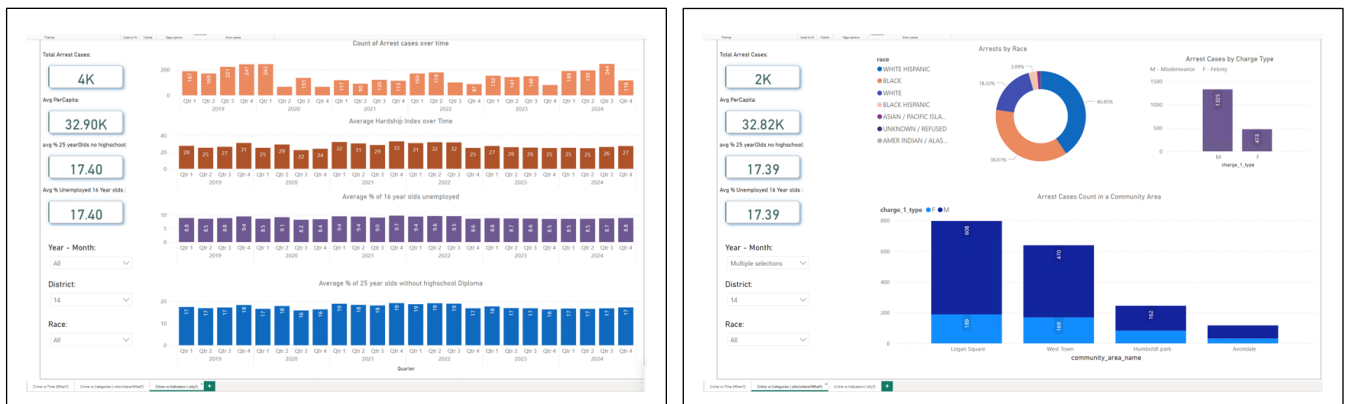
## 4. Serving/Analysis

The joined data is stored back into Azure Blob Storage, specifically within the cleaned blob container. From there, the data is imported into Power BI for analysis. Additionally, a copy of the aggregated datasets is loaded into Jupyter Notebook for further analysis, as outlined in the overall technical architecture.

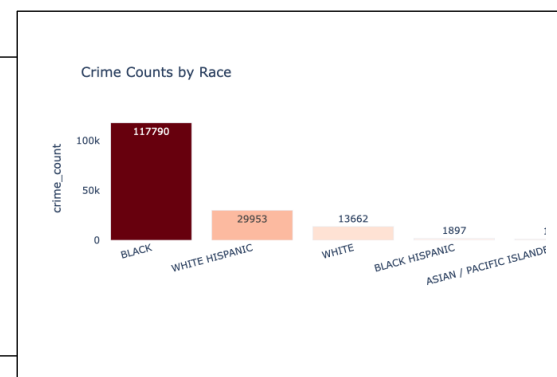
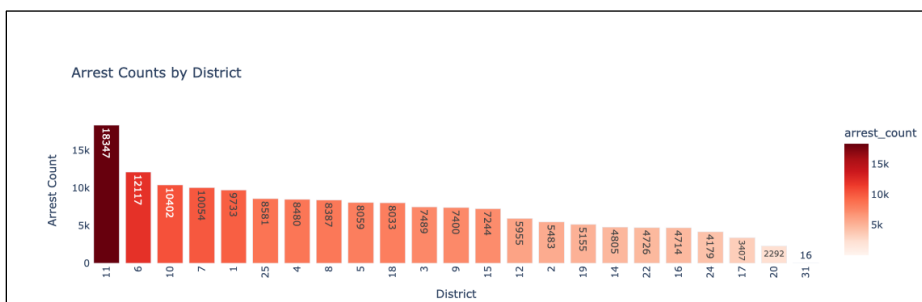
#### 4.1 Guiding Questions for Analysis:

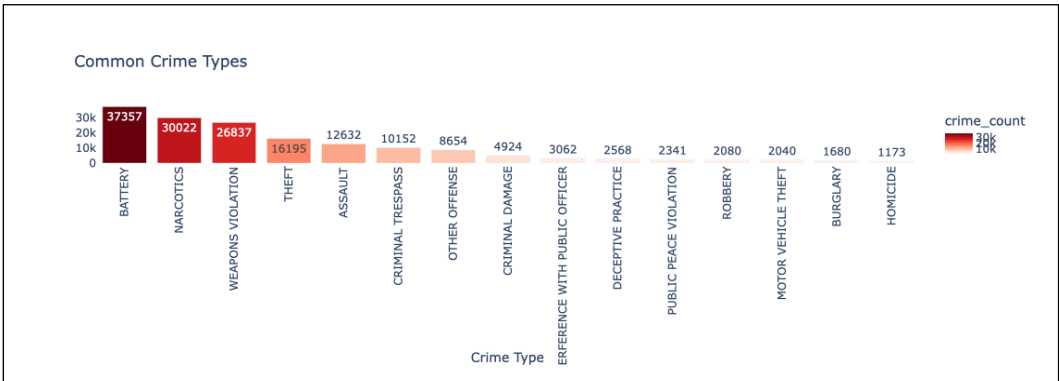
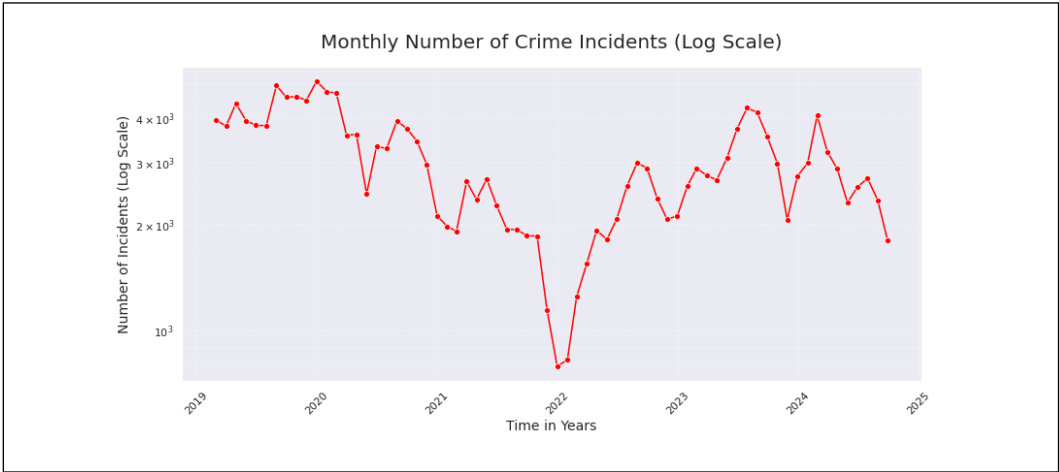
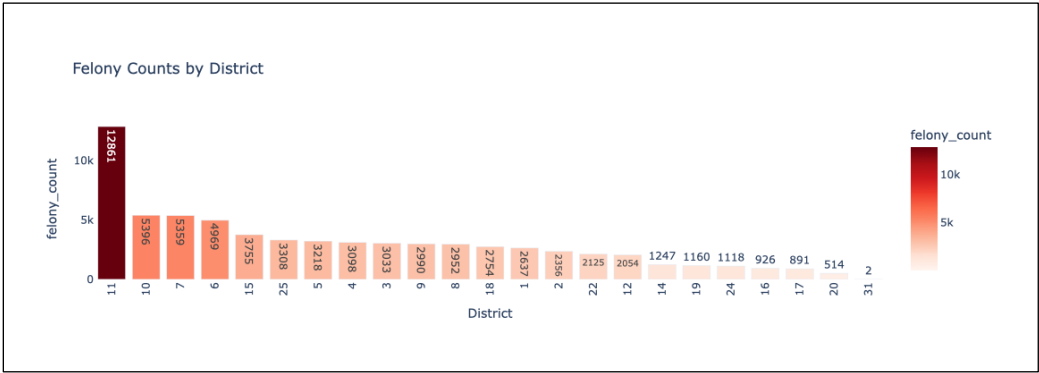
- What areas with high crime rates relative to socioeconomic factors?
- How crime trended over years, Months?
- Is crime more common towards weekends or weekdays?
- What are High-risk districts that need more attention?
- Does poor socio economic indicators imply high crime rate?

#### 4.2 Visuals from Power BI:

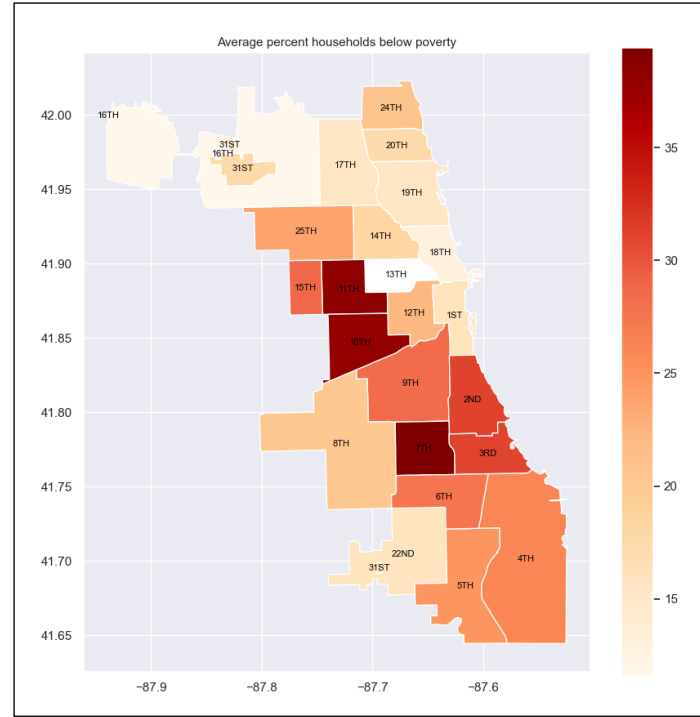
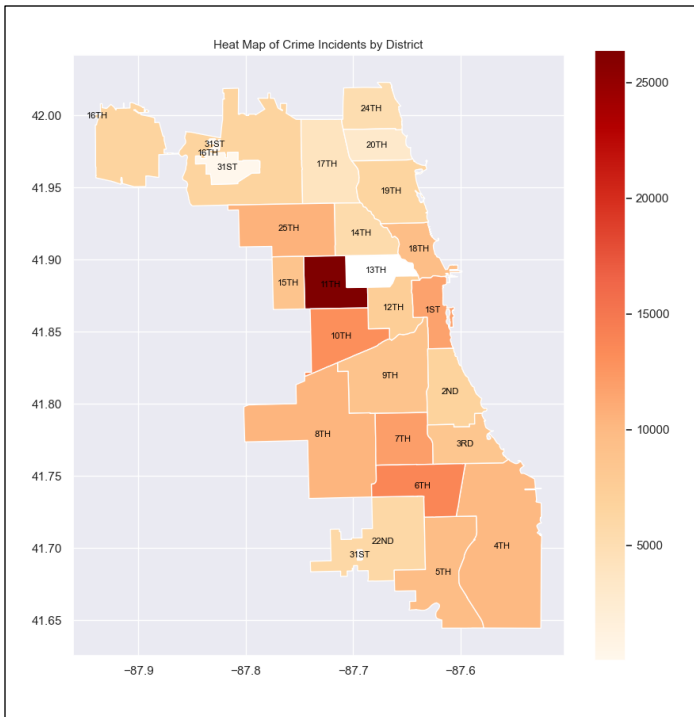
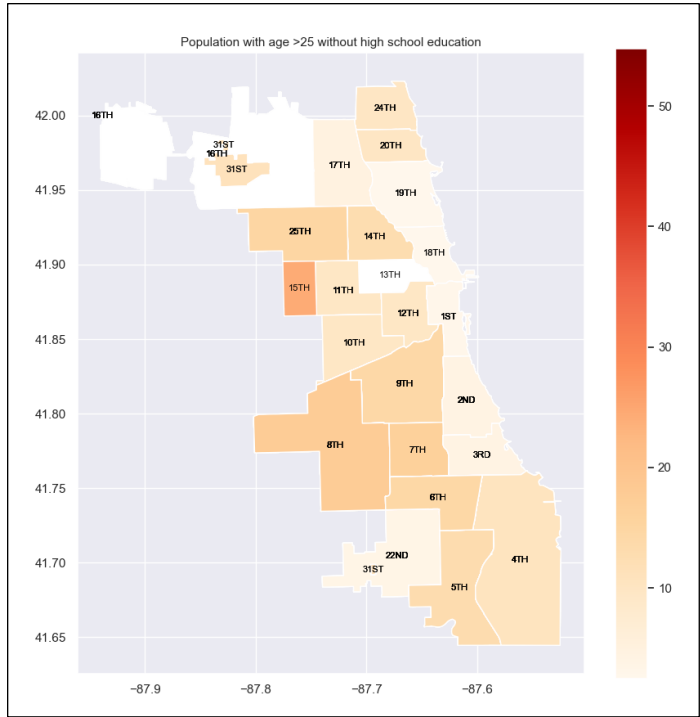
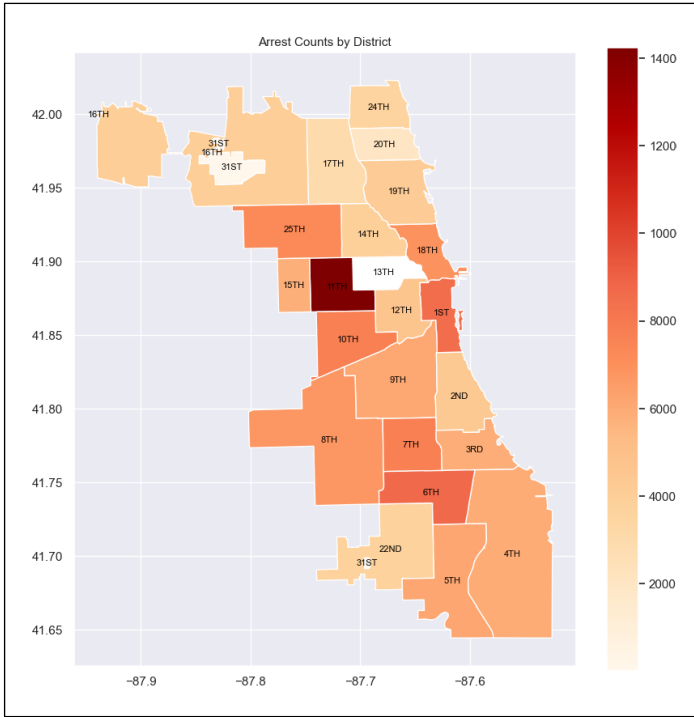


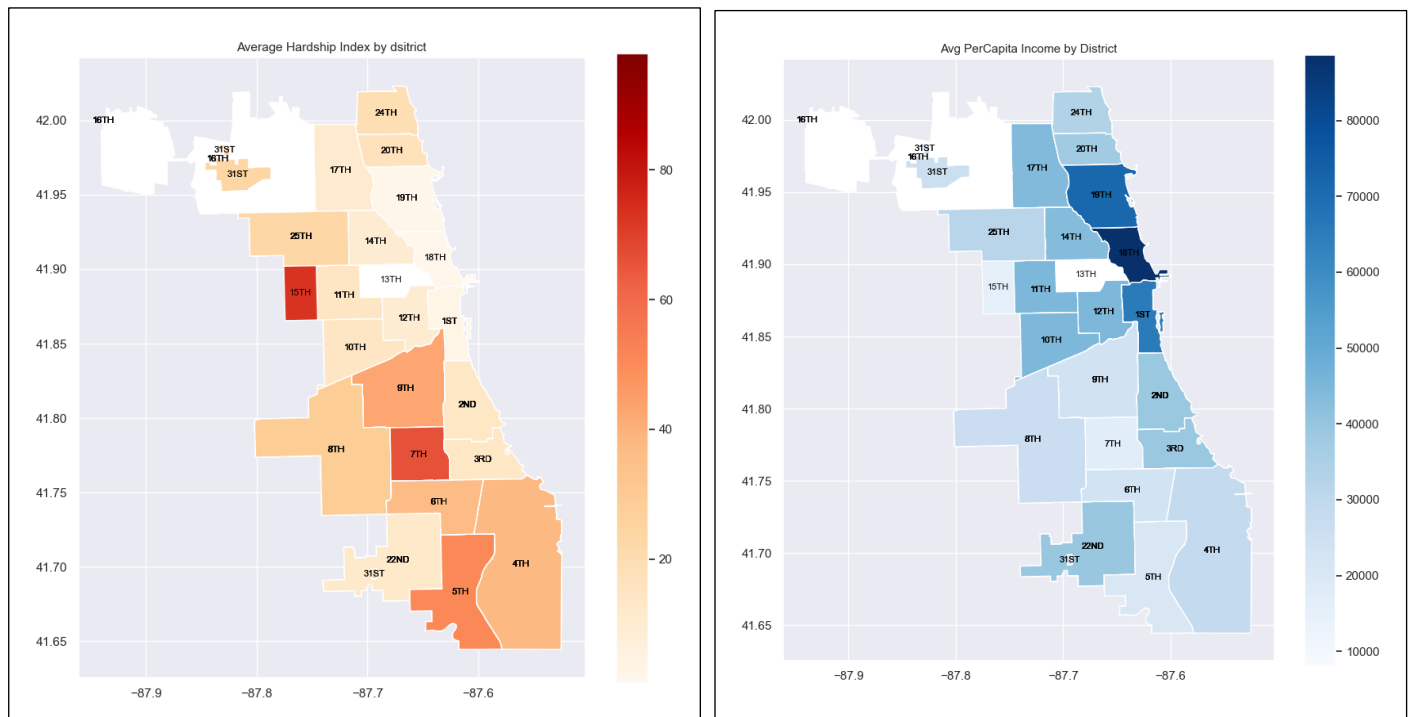
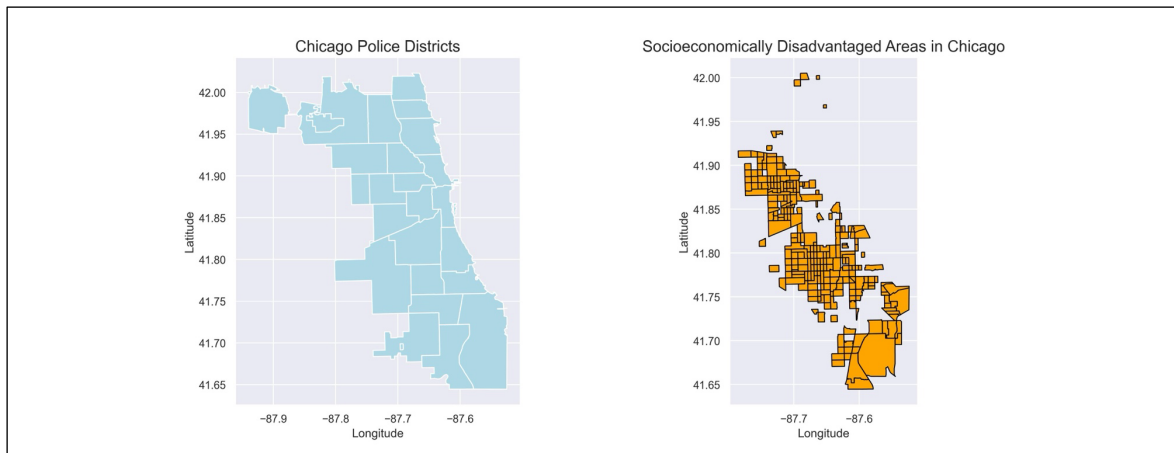
#### 4.3 Visuals from Jupyter Notebook











#### 4.4 Analysis Insights:

- Police District 11 in Chicago has highest rate of arrests.
- Most of the crime incidents happen on the street.
- Crime rates tend to rise towards weekend days of Friday, Saturday and Sunday.
- Crime rates seem higher towards march and April based on analysis of 5-year aggregated data.
- The most common crime types are Battery, narcotics, weapons violations and theft.
- Socio economic indicators seem to have a positive correlation with crime as socio-economically disadvantaged areas experience more crime.

#### 4.5 Actions for call based on Analysis Insights:

- Pay more attention or allocate resources for specific police districts such as District 11 or on weekends.
- To prevent most crime incidents that happen on the street, police must act.
- Pay attention to most common types of crimes such as Battery, narcotics, weapons violations and theft.
- Pay attention to Socio economically disadvantaged areas to increase Per-Capita income and education level to proactively reduce crime rates.

#### 5. Further objectives to improve the project

- Automate ingestion with automation tools like Azure data factory or Airflow.
- Automatic deployment of cloud using Infrastructure as a code (Terraform)
- Unit testing of code modules