# Individual Assignment

Je Sai Kailash Pulipati

002339774

## Overview

*Basic Information*

- Total Number of records/rows/lines: 260,503 records
- Total Number of columns/fields: 19 Columns
- Source and Target Data types
- Grain: Each individual arrest record

| Column Name | Source Data Type | Target Data Type (SQL) |
|---|---|---|
| ARREST_KEY | String | VARCHAR |
| ARREST_DATE | V_String | DATE |
| PD_CD | V_String | INT |
| PD_DESC | V_String | VARCHAR |
| KY_CD | String | INT |
| OFNS_DESC | String | VARCHAR |
| LAW_CODE | String | VARCHAR |
| LAW_CAT_CD | String | VARCHAR |
| ARREST_BORO | V_String | VARCHAR |
| ARREST_PRECINCT | String | INT |
| JURISDICTION_CODE | V_String | INT |
| AGE_GROUP | V_String | VARCHAR |
| PERP_SEX | String | VARCHAR |
| PERP_RACE | String | VARCHAR |
| X_COORD_CD | V_String | INT |
| Y_COORD_CD | V_String | INT |
| Latitude | V_String | FLOAT |
| Longitude | V_String | FLOAT |

| New Georeferenced Column | String | VARCHAR |
|---|---|---|

## Data Quality Assessment and Data Cleaning

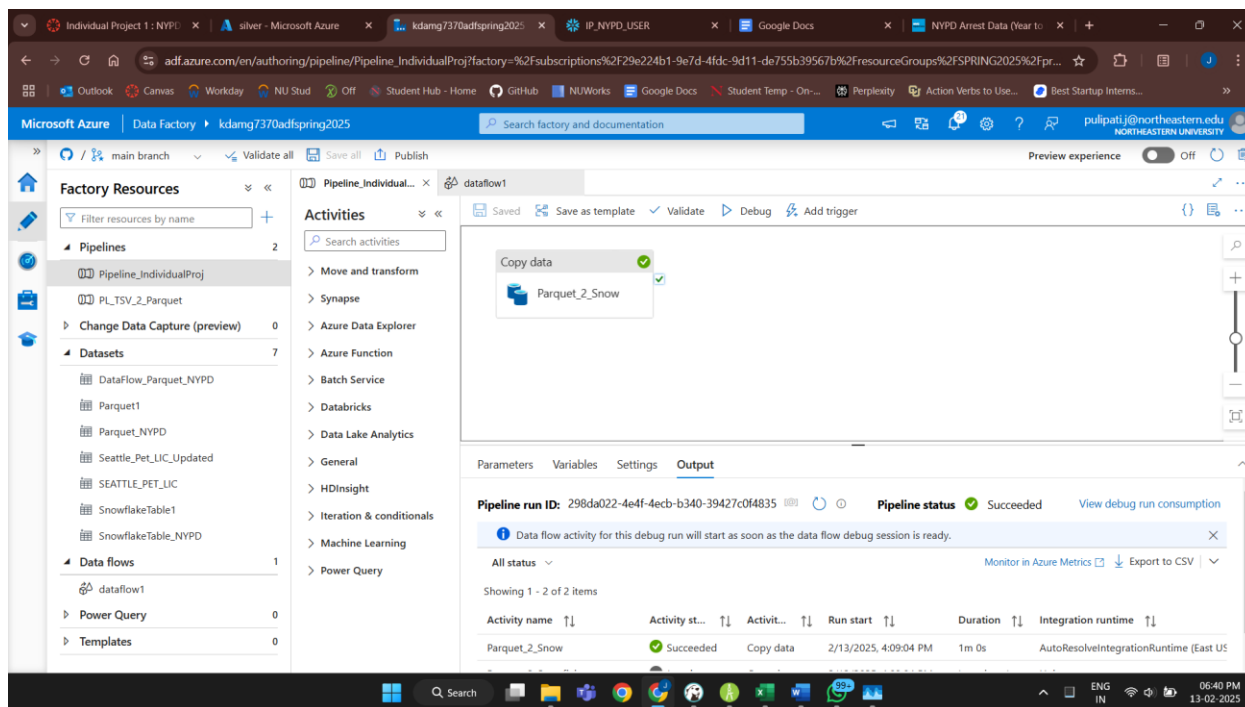| Sno. | Column Name | Issue | Details | Action Plan to fix it |
|---|---|---|---|---|
| 1. | ARREST_DATE | Date in String format | Need to convert it to date format | Use parse DateTime to change the format to DD-MM-YYYY |
| 2. | LAW_CAT_CD | Only 3 levels of Felony present | Other types or Misc felonies reported like "9", "I", "(null)" | Used filter Tool to remove unwanted rows |
| 3. | ARREST_DATE, PD_CD, KY_CD, | Null Values present | Contains Null values in between | Used Data cleaning tool to remove the null values and replace with "0" |
| 4. | PD_DESC | Null Values present | Contains Null values like "(null)" | Used formula tool to replace "(null)" to "Other" |
| 5. | Longitude and Latitude | Null Values | Contains Null values | Used Data Imputation to replace null values with Mean |
| 6. | Select Tool | Wrong Data format | Need to change Data formats | Used Select tool to change the data types |
| 7. | AGE_GROUP | Unusable age group divisions | Need to remove divisions and use some unique identifying values | Used Formula tool to divide age groups into values |
| 8. | LAW_CAT_CD | Jurisdiction out of range | Only 0,1,2 Jurisdictions should be present | Use filter tool to separate rows which are above 2 |

## Insights and Observations

1. The LAW_CAT_CD column has violations other than M, F, V.
2. There are some Jurisdiction codes which are outside the NYPD which can be removed
3. Null values are identified using the Summarise tool
4. There are no duplicate Arrest_key records
5. Some longitudes, latitudes have "0" so used data imputation to fill with average values so that they won't be affected by becoming outliers
6. Converted Age_groups into categories like 1,2,3,4,5.
7. Some column names are not in proper stanndart format, so converted all column names into Capital with no spaces.
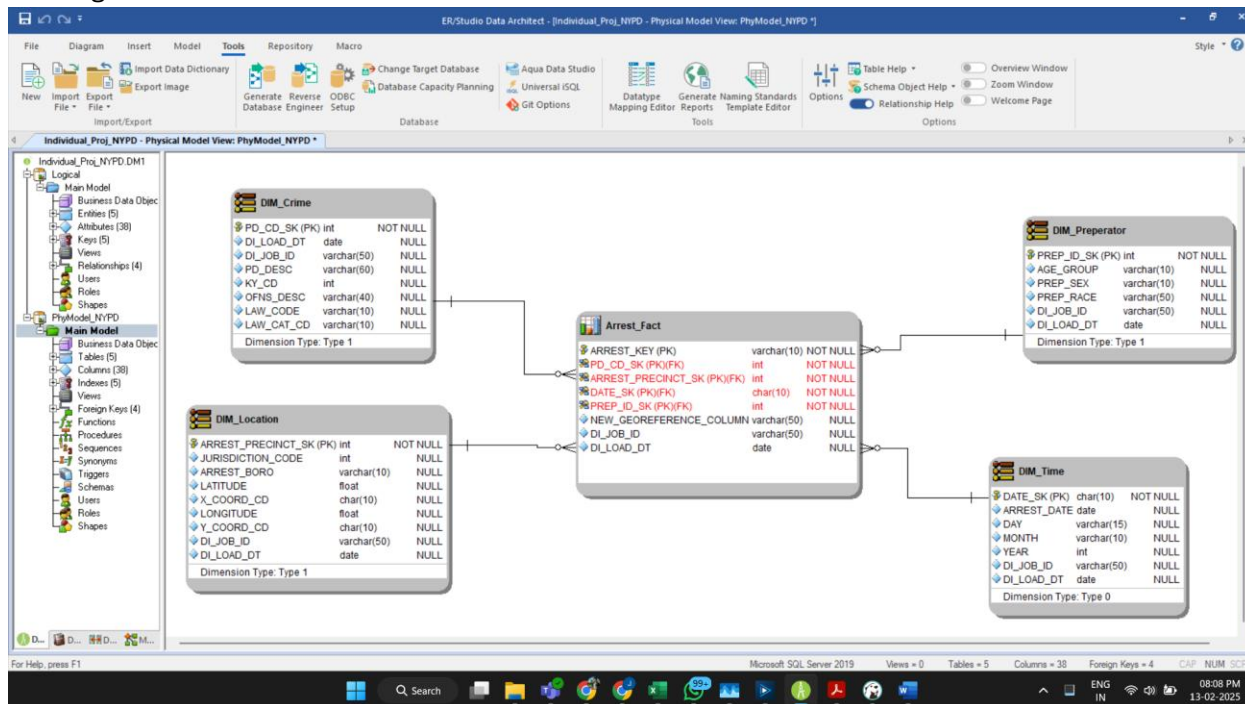
## Alteryx Workflow



## ADF-Snowflake Pipeline

ER Diagram:

**Final Created Table pushed to Snowflake:**



**Snowflake DB Creation**

# Count (*) rows from Table



# DIMENSIONS AND FACTS TABLES:

SQL SCRIPT

/*

 * ER/Studio Data Architect SQL Code Generation

 * Project :     Individual_Proj_NYPD.DM1

 *

 * Date Created : Thursday, February 13, 2025 20:06:58

 * Target DBMS : Microsoft SQL Server 2019

 */


/*

 * TABLE: Arrest_Fact

 */


CREATE TABLE Arrest_Fact(

    ARREST_KEY            varchar(10)   NOT NULL,

    PD_CD_SK             int        NOT NULL,

    ARREST_PRECINCT_SK       int        NOT NULL,

```sql
    DATE_SK              char(10)    NOT NULL,

    PREP_ID_SK           int         NOT NULL,

    NEW_GEOREFERENCE_COLUMN   varchar(50)   NULL,

    DI_JOB_ID            varchar(50)   NULL,

    DI_LOAD_DT           date          NULL,

    CONSTRAINT PK_ARREST_KEY PRIMARY KEY NONCLUSTERED (ARREST_KEY, PD_CD_SK,
ARREST_PRECINCT_SK, DATE_SK, PREP_ID_SK)

)


go



IF OBJECT_ID('Arrest_Fact') IS NOT NULL

    PRINT '<<< CREATED TABLE Arrest_Fact >>>'

ELSE

    PRINT '<<< FAILED CREATING TABLE Arrest_Fact >>>'

go


/*

 * TABLE: DIM_Crime

 */


CREATE TABLE DIM_Crime(

    PD_CD_SK    int        NOT NULL,

    DI_LOAD_DT  date        NULL,

    DI_JOB_ID   varchar(50)  NULL,

    PD_DESC     varchar(60)  NULL,

    KY_CD       int          NULL,

    OFNS_DESC   varchar(40)  NULL,
```

```sql
    LAW_CODE    varchar(10)   NULL,

    LAW_CAT_CD   varchar(10)   NULL,

    CONSTRAINT PK2 PRIMARY KEY NONCLUSTERED (PD_CD_SK)
)


go



IF OBJECT_ID('DIM_Crime') IS NOT NULL

    PRINT '<<< CREATED TABLE DIM_Crime >>>'
ELSE

    PRINT '<<< FAILED CREATING TABLE DIM_Crime >>>'
go


/*
 * TABLE: DIM_Location
 */


CREATE TABLE DIM_Location(

    ARREST_PRECINCT_SK   int       NOT NULL,

    JURISDICTION_CODE    int        NULL,

    ARREST_BORO        varchar(10)   NULL,

    LATITUDE          float       NULL,

    X_COORD_CD        char(10)     NULL,

    LONGITUDE         float       NULL,

    Y_COORD_CD        char(10)     NULL,

    DI_JOB_ID        varchar(50)   NULL,

    DI_LOAD_DT        date        NULL,

    CONSTRAINT PK_ARREST_PRECINCT PRIMARY KEY NONCLUSTERED (ARREST_PRECINCT_SK)
```

```sql
)

go


IF OBJECT_ID('DIM_Location') IS NOT NULL
   PRINT '<<< CREATED TABLE DIM_Location >>>'
ELSE
   PRINT '<<< FAILED CREATING TABLE DIM_Location >>>'
go


/*
 * TABLE: DIM_Preperator
 */


CREATE TABLE DIM_Preperator(
   PREP_ID_SK   int        NOT NULL,
   AGE_GROUP    varchar(10)   NULL,
   PREP_SEX    varchar(10)   NULL,
   PREP_RACE    varchar(50)   NULL,
   DI_JOB_ID    varchar(50)   NULL,
   DI_LOAD_DT   date        NULL,
   CONSTRAINT PK4 PRIMARY KEY NONCLUSTERED (PREP_ID_SK)
)


go



IF OBJECT_ID('DIM_Preperator') IS NOT NULL
```

```
    PRINT '<<< CREATED TABLE DIM_Preperator >>>'
ELSE
    PRINT '<<< FAILED CREATING TABLE DIM_Preperator >>>'
go


/*
 * TABLE: DIM_Time
 */


CREATE TABLE DIM_Time(
    DATE_SK      char(10)    NOT NULL,
    ARREST_DATE  date        NULL,
    DAY          varchar(15)  NULL,
    MONTH        varchar(10)  NULL,
    YEAR         int         NULL,
    DI_JOB_ID    varchar(50)  NULL,
    DI_LOAD_DT   date         NULL,
    CONSTRAINT PK_DATE_ID PRIMARY KEY NONCLUSTERED (DATE_SK)
)

go



IF OBJECT_ID('DIM_Time') IS NOT NULL
    PRINT '<<< CREATED TABLE DIM_Time >>>'
ELSE
    PRINT '<<< FAILED CREATING TABLE DIM_Time >>>'
go
```

```
/*
 * TABLE: Arrest_Fact
 */


ALTER TABLE Arrest_Fact ADD CONSTRAINT RefDIM_Crime1
    FOREIGN KEY (PD_CD_SK)
    REFERENCES DIM_Crime(PD_CD_SK)
go


ALTER TABLE Arrest_Fact ADD CONSTRAINT RefDIM_Location2
    FOREIGN KEY (ARREST_PRECINCT_SK)
    REFERENCES DIM_Location(ARREST_PRECINCT_SK)
go


ALTER TABLE Arrest_Fact ADD CONSTRAINT RefDIM_Time3
    FOREIGN KEY (DATE_SK)
    REFERENCES DIM_Time(DATE_SK)
go


ALTER TABLE Arrest_Fact ADD CONSTRAINT RefDIM_Preperator4
    FOREIGN KEY (PREP_ID_SK)
    REFERENCES DIM_Preperator(PREP_ID_SK)
go
```

**DDL Scripts**

1. How many arrests occurred on any specific day, week, month, quarter, or year?

SELECT dt.YEAR, dt.MONTH, COUNT(*) AS Total_Arrests FROM Arrest_Fact af

JOIN DIM_Time dt ON af.DATE_SK = dt.DATE_SK GROUP BY dt.YEAR, dt.MONTH

ORDER BY dt.YEAR, dt.MONTH;


2. What are the peak days and months for arrests?

SELECT dt.MONTH, COUNT(*) AS Arrest_Count FROM Arrest_Fact af

JOIN DIM_Time dt ON af.DATE_SK = dt.DATE_SK GROUP BY dt.MONTH ORDER BY Arrest_Count DESC;


3. What are the top 5 most frequently occurring crimes?

SELECT dc.OFNS_DESC, COUNT(*) AS Crime_Count

FROM Arrest_Fact af JOIN DIM_Crime dc ON af.PD_CD_SK = dc.PD_CD_SK

GROUP BY dc.OFNS_DESC ORDER BY Crime_Count DESC;


4. Which crimes have increased or decreased the most over time?

SELECT dt.YEAR, dc.OFNS_DESC, COUNT(*) AS Arrest_Count FROM Arrest_Fact af JOIN DIM_Time dt ON af.DATE_SK = dt.DATE_SK JOIN DIM_Crime dc ON af.PD_CD_SK = dc.PD_CD_SK GROUP BY dt.YEAR, dc.OFNS_DESC ORDER BY dt.YEAR, Arrest_Count DESC;


5. Are there specific precincts with higher felony arrests compared to misdemeanors? (Hint: A precinct is a police district within a city.)

SELECT dl.ARREST_PRECINCT, dc.LAW_CAT_CD, COUNT(*) AS Arrest_Count

FROM Arrest_Fact af JOIN DIM_Location dl ON af.ARREST_PRECINCT_SK = dl.ARREST_PRECINCT_SK

JOIN DIM_Crime dc ON af.PD_CD_SK = dc.PD_CD_SK

WHERE dc.LAW_CAT_CD IN ('F', 'M') GROUP BY dl.ARREST_PRECINCT, dc.LAW_CAT_CD ORDER BY Arrest_Count;


6. Which borough has the highest number of arrests? (Hint: A borough is a large administrative division in NYC, such as Manhattan (M), Brooklyn (K), Queens (Q), The Bronx (B), and Staten Island (S).)

SELECT dl.ARREST_BORO, COUNT(*) AS Arrest_Count FROM Arrest_Fact af JOIN DIM_Location dl ON af.ARREST_PRECINCT_SK = dl.ARREST_PRECINCT_SK

GROUP BY dl.ARREST_BORO ORDER BY Arrest_Count DESC;

7.  What is the distribution of arrestees by age, race, and gender?

SELECT dp.AGE_GROUP, dp.PREP_RACE, dp.PREP_SEX, COUNT(*) AS Arrest_Count

FROM Arrest_Fact af JOIN DIM_Preperator dp ON af.PREP_ID_SK = dp.PREP_ID_SK

GROUP BY dp.AGE_GROUP, dp.PREP_RACE, dp.PREP_SEX ORDER BY Arrest_Count DESC;

8.  Can we predict high-crime areas based on past arrest data?

SELECT dl.LATITUDE, dl.LONGITUDE, COUNT(*) AS Arrest_Count FROM Arrest_Fact af JOIN DIM_Location dl ON af.ARREST_PRECINCT_SK = dl.ARREST_PRECINCT_SK GROUP BY dl.LATITUDE, dl.LONGITUDE ORDER BY Arrest_Count DESC;