

GROUP 2

Crimes to Arrests

COURSE NUMBER: CIS 4400

COURSE SECTION: CMWA

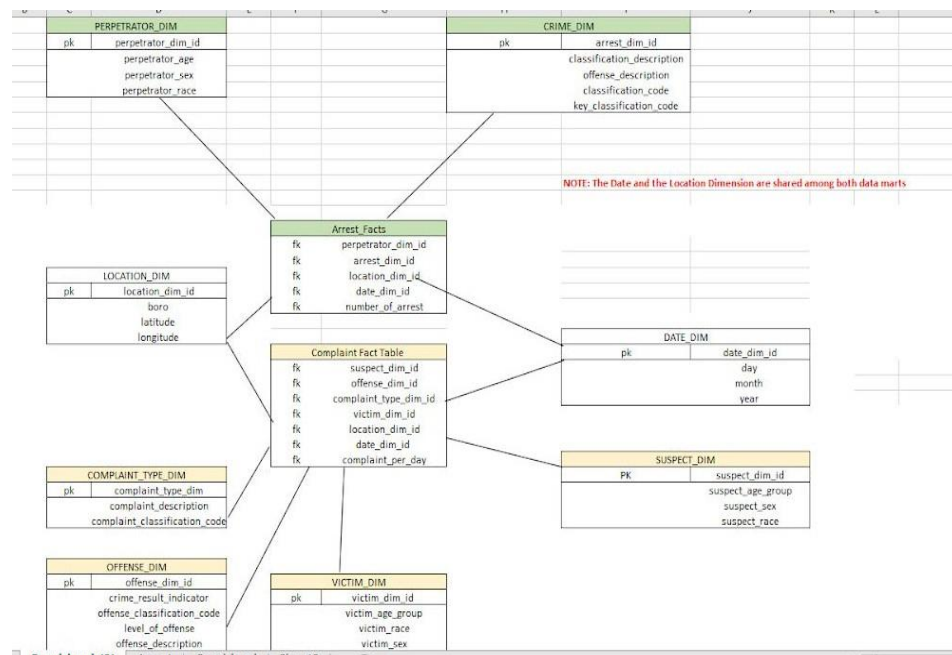
1. An introduction section similar to the proposal section including the narrative description of the business used for the data warehouse being created and description of the source data.

The business that we are conducting is a non-profit that is trying to see the most affected population who is arrested. We are trying to see if there are any attributes that contribute the most to the arrested population and tracking if all the crimes reported lead to arrests.

The data sources that were used for this project are derived from NYC Open Data. We found the NYPD Complaint Data set which contains all valid felonies, misdemeanors, and violation crimes reported to the NYPD; this data set includes all records beginning from 2016 and the granularity of this data is quarterly. This data set contains 35 columns but for the purpose of this project we only chose the important attributes which you will see in the dimensional model below.

Next up we have the NYPD Arrest Data set which is a breakdown of every arrest in NYC by the NYPD. Like the data set above, it is updated quarterly but contains those records beginning from 2018. The set has 18 columns but did not use all of them as shown below by our dimensional diagram.

2. Dimensional Model Diagram

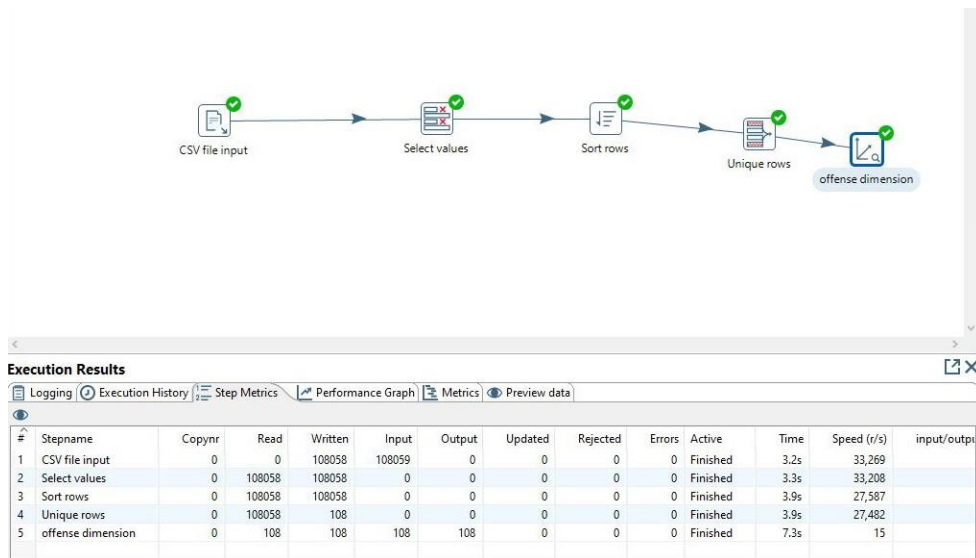


3. A description and screen pictures of the ETL process

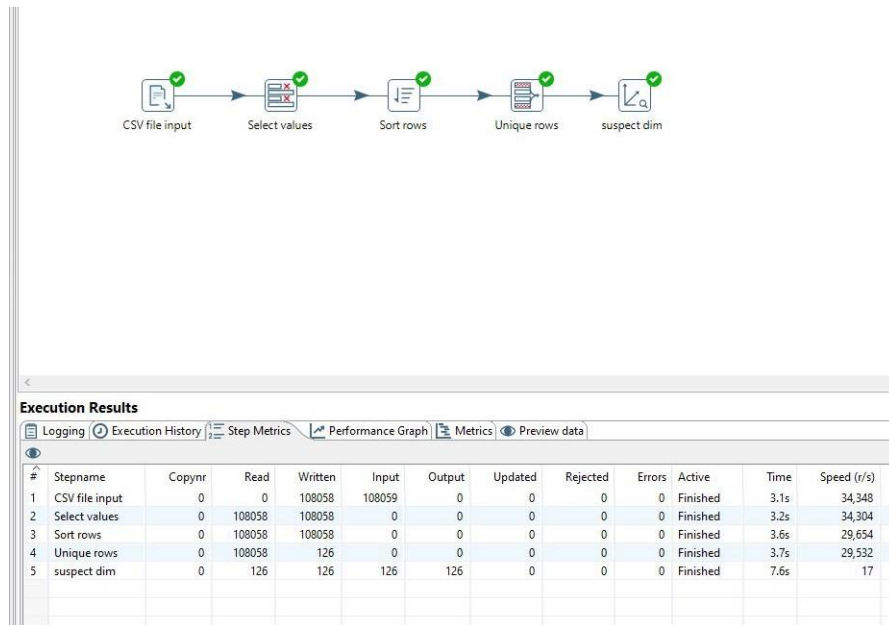
Description: This transformation started off with the CSV file. Then we selected off the values for the dimension, sorted and picked out the unique values from the set. Finally we created the **complaint** dimension



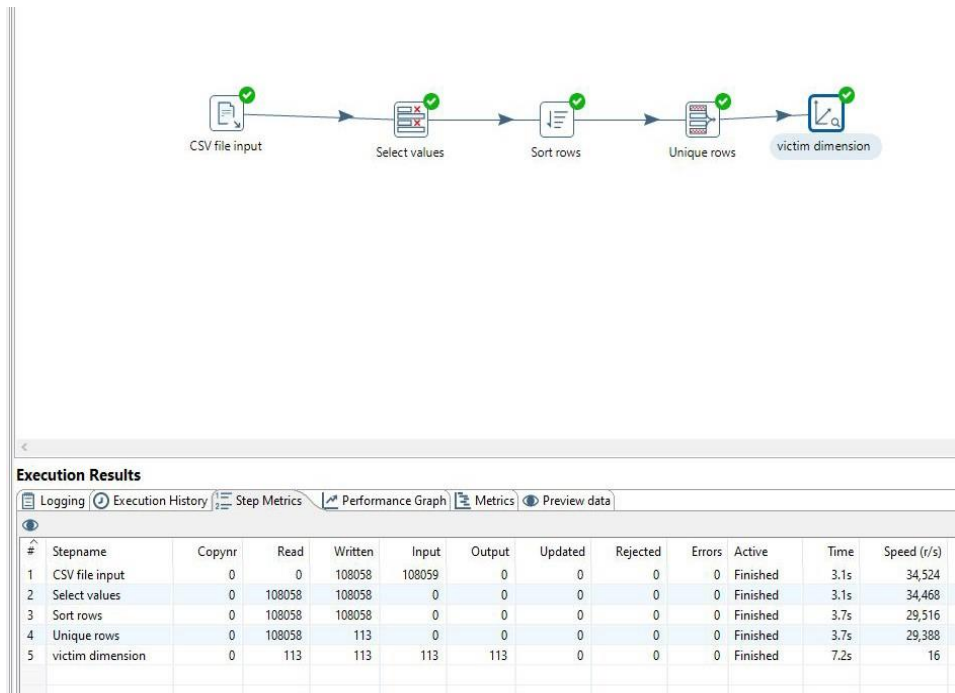
Description: Following the same steps as before we created the **offense** dimension which contained 4 of the original attributes from the dataset.



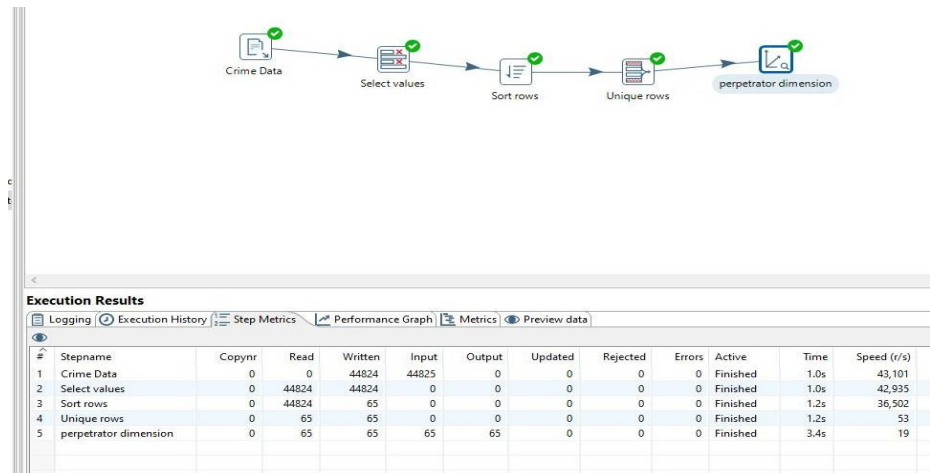
Description: This is the **suspect** dimension which contains the suspects age group, sex, and race.



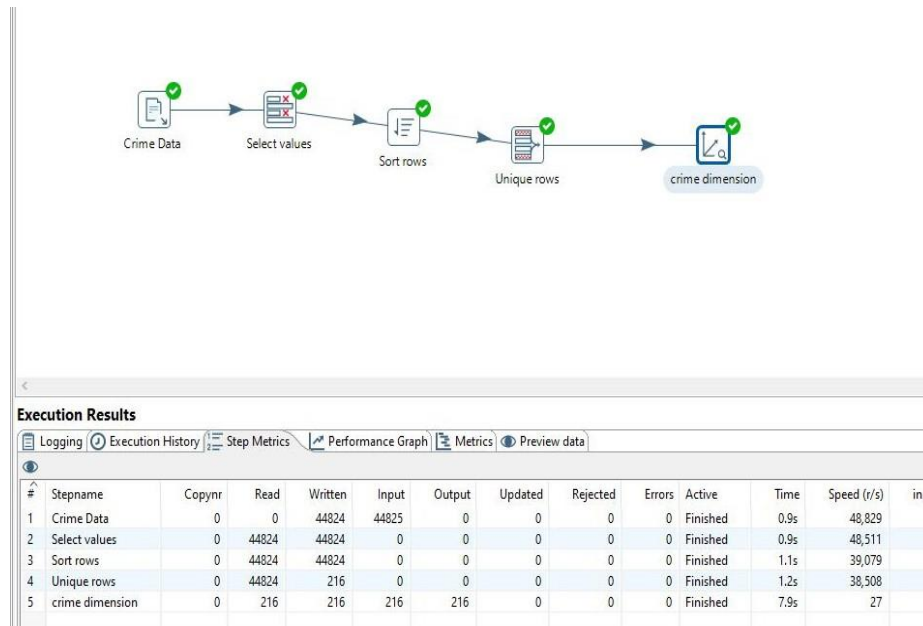
Description: This is the **victim** dimension and similar to the one above this one contains the victims age group, sex, and race



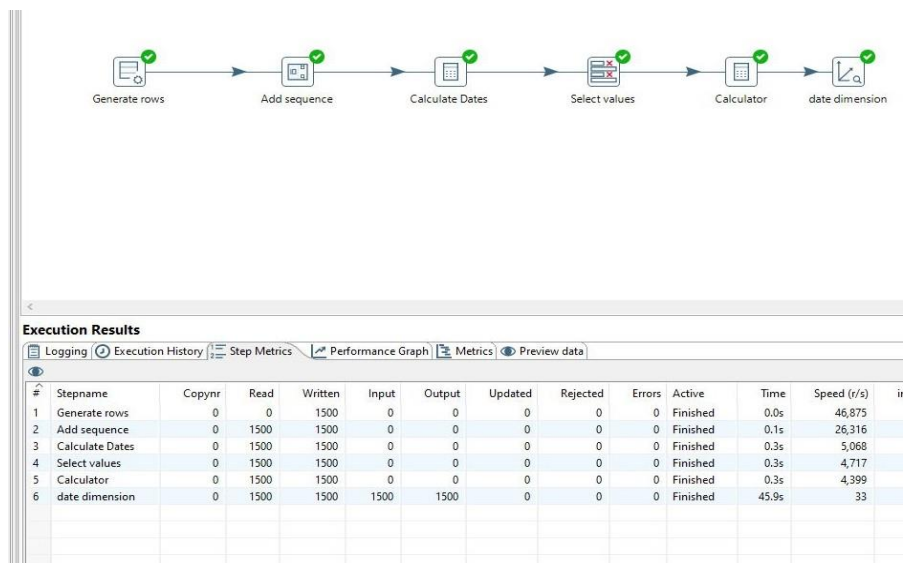
Description: This is the **perpetrator** dimension and similar to the one above it contains the perpetrators age group, sex, and race



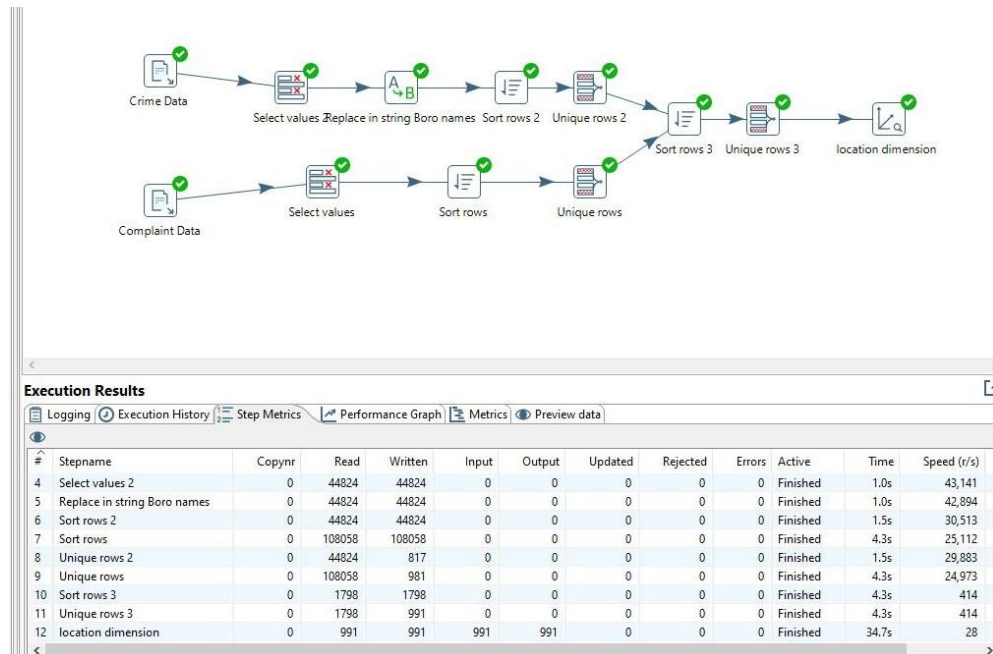
Description: This is the **crime** dimension which contains original attributes from the data set including the classification code, classification description, key classification code and the offense description.



Date and Location Dimensions

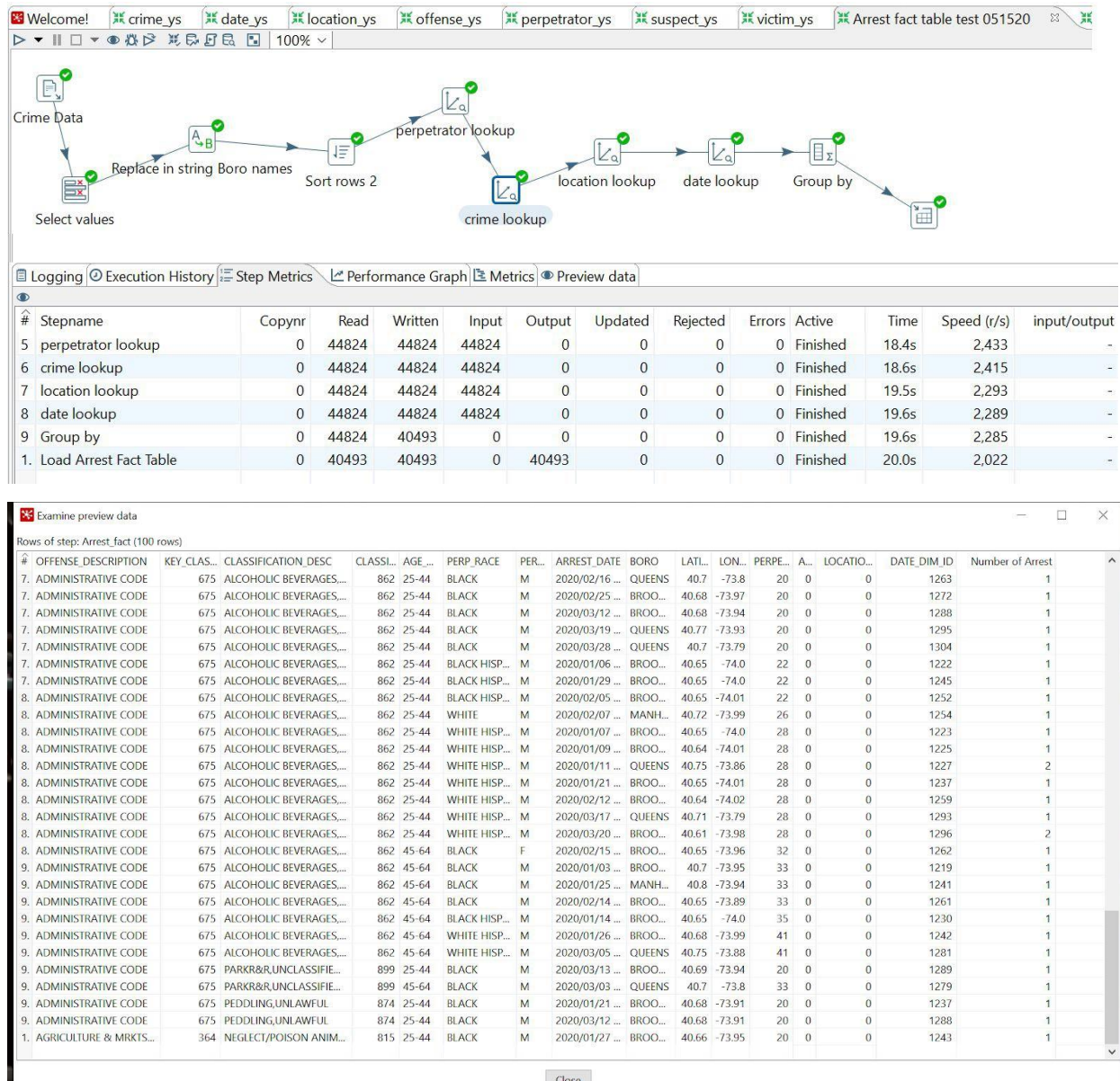


Description: This is the date dimension for both sets of data. Because one of them starts off in 2016 and the other in 2018 we needed to make sure that they have enough dates to measure up to the last time both sets were updated. We calculated the dates so that it shows up for four years. Also, since our model wants only the date, month, and year we did another calculated step where we could see the month and year only



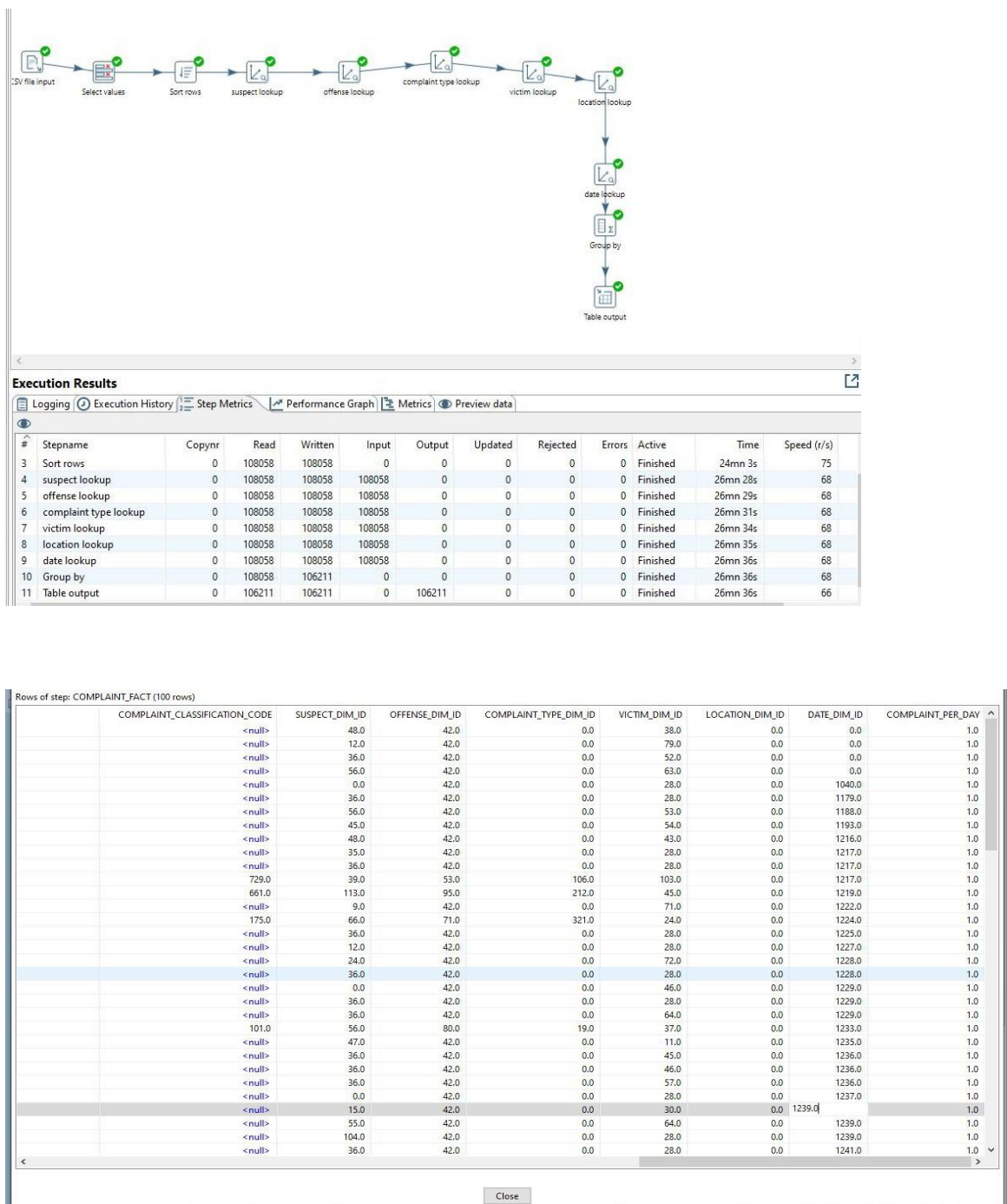
Description: This is the location dimension and this is what will bind together the entire data set in addition to the date dimension. We took the CSV files from both data sets and filtered through them to get the values for location (borough name, longitude, latitude). For one of the inputs, the borough name was limited to one character so we needed to incorporate the 'replace with string' step to make sure that they both aligned with each other. We also rounded two decimal points for this dimension. After filtering and sorting through the records we ended up with 991 unique locations.

4. Screen shots and brief descriptions of the final schema that the business analytics tools are working with.

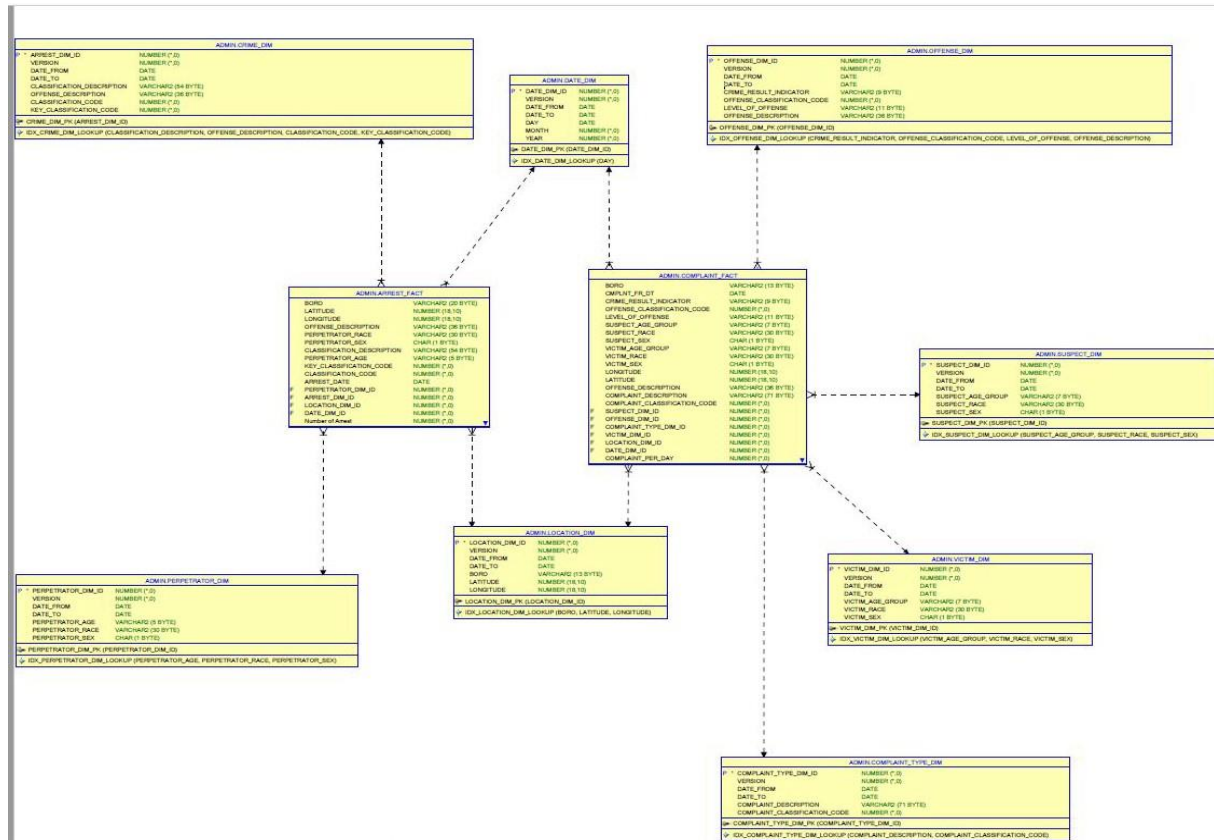


Arrest Fact Table: Reads a CSV file from NYPD_Arrest_Data_Year_to_Date (2020 from January 1 to March 31, 2020) data file. The total number of rows turns out to be 44824. “Group By” on Pentaho enables a column to be duplicated but counted by the number of N. It will represent “1” per perpetrator_dim_number, so that it can be calculated for KPI (arrests per day, arrests per age group and etc).

Complaint Fact Table:

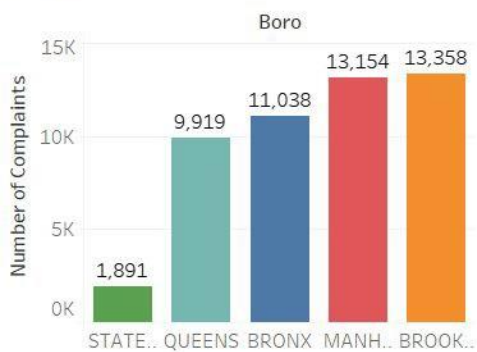


Description: This is the final schema that we reverse engineered in Oracle SQL Developer. The relationships were made by altering the table to create the foreign key relationships.

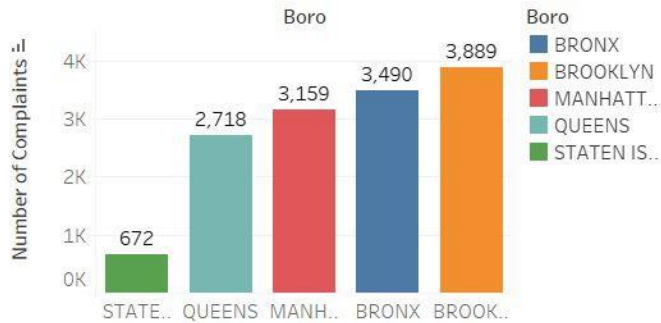


5. Screen shots and descriptions of the analytics (at least 3) on the dashboard application developed based on the data warehouse data.

Number of Complaints per Male



Number of Complaints per Female



Number of Male Arrests per Age Group

Perpe	trator
Age	#
<18	1,474
18-24	6,419
25-44	18,179
45-64	6,864
65+	464

Number of Female Arrests per Age Group

Perpe	trator
Age	#
<18	341
18-24	1,471
25-44	3,970
45-64	1,205
65+	73

Description: As most of our KPIs were based on qualitative values we decided to stick with bar graphs and basic charts to reflect the number that we were looking for in our KPI's. These KPIs tell us a lot about the data set. The first being that the complaints per male are larger than that of females. We can also tell that both female and male complaints take place in the borough of Brooklyn more than any other borough. When looking at the arrests, we can see that there are more male arrests than female arrests which makes sense as they received more complaints than females. Looking from the basic count of arrests we can determine that the age group most likely to be arrested is from 25-44 and from this we can infer the same for the number or complaints per age group. Below are the KPIs we looked at and the code for each.

Number of Complaints per Male

```
SELECT complaint_description, suspect_sex, boro
FROM complaint_fact
WHERE complaint_description IS NOT NULL AND boro IS NOT NULL AND suspect_sex =
'M';
```

Number of complaints per female

```
SELECT complaint_description, suspect_sex, boro
FROM complaint_fact
WHERE complaint_description IS NOT NULL AND boro IS NOT NULL AND suspect_sex =
'F';
```

Number of arrests per age group and are female

```
SELECT offense_description, perpetrator_age, boro, perpetrator_sex
FROM arrest_fact
WHERE offense_description IS NOT NULL AND boro is NOT NULL AND perpetrator_sex
= 'F'
```

Number of arrests per age group and are male

```
SELECT offense_description, perpetrator_age, boro, perpetrator_sex
FROM arrest_fact
WHERE offense_description IS NOT NULL AND boro is NOT NULL AND perpetrator_sex
= 'M'
```

6. A narrative conclusion section that describes:

After concluding our project there were many challenges that came along the way. The first one that we came across was loading all of the data from each dimension. After figuring out which schema to finally put everything on, we decided to work on the fact tables. It was difficult as we kept getting errors everytime we ran the fact tables which prevented us from getting to the next step. The errors mostly came from our Oracle cloud schema since it does not allow multiple users at the sametime when using the free version. However, we were able to eventually figure out that by separating the work on localhost (H2 database) and sharing the kettle files we would be able to implement working on dimensions or fact tables into the Oracle database so that everyone could stay and catch up to the same step of the project. The new system shows users the number of complaints or arrests per day across New York City, and can also be observed by their location /area. It will contribute to the police and fire department in looking at

the big picture to increase their operational efficiency in which areas they should place focus their efforts on and place more of their staff on.