



Final Project
DAV-6100-Information Architecture
Group 3



Yeshiva University | KATZ SCHOOL

NYC Car Crash Data Warehouse Design

Overview

- Our project centers on analyzing and interpreting the Motor Vehicle Collisions data for New York City.
- The DW address critical questions surrounding incident patterns, location, collusion trends over time, and factors contributing to collisions.
- We aim to uncover insights that could bolster traffic management strategies, enhance public safety initiatives, and inform policy decisions.
- The project aims also at facilitating a better understanding of traffic dynamics in NYC, which could lead to improved road safety measures and a reduction in collision rates.
- The anticipated outcome is a robust data-driven framework that provides NYC authorities, policymakers, and the public with the tools necessary to comprehend and improve the multifaceted nature of urban traffic safety.



Roles & Responsibilities

Rekha & Jetendra - Data Engineer:

Integrates and manages data workflows into AWS cloud architecture. Also focuses on efficiency and reliability of data collection and processing.

Priyank & Volkan - Data Analyst:

Interprets data patterns, providing insights and recommendations to support strategic decision-making.



Timeline of the Project

Week 4: Bus Matrix, Basic Use Cases

Week 9: Architecture Diagram

Week 5: Dimensional Model

Week 10-12: ETL/ ELT

Week 6-7: Conceptual, Logical and Physical Models

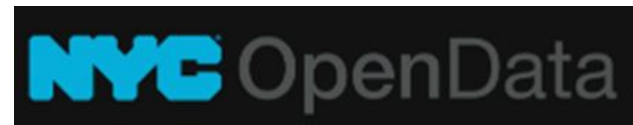
Week 13-14: Analytical Queries and Visualization

Week 8: Data Dictionary



About Data

Motor Vehicle Collisions in NYC - records comprehensive details from all police-reported motor vehicle collisions in New York City. It requires the completion of police reports for incidents involving injuries, fatalities, or property damage exceeding \$1000. It records various aspects such as the specifics of each vehicle involved, personal details of individuals affected (drivers, passengers, pedestrians, etc.), and the overall crash event characteristics. This data, which dates back to the adoption of the electronic reporting system in April 2016, serves as a pivotal resource for analyzing collision causes, assessing traffic safety policies, and developing preventative measures.



Data Source:

<https://data.cityofnewyork.us/Public-Safety/Motor-Vehicle-Collisions-Vehicles/bm4k-52h4>



Summary of Datasets

Entities

Date: Represents the date and time aspects of a crash.

Contribution Factor:
Captures the factors contributing to a crash.

Vehicle: Details about the vehicles involved in the crashes.

Person: Information about individuals involved in the crashes.

Crash Event: The central entity that records each crash incident.

Method of access

- Pull data from internet open-source
- MySQL connector and Python (Jupyter Notebook)
- Load into AWS S3 bucket

Data Quality

- The dataset includes fair amount of missing or incomplete information.
- We Identified one dataset is not been mostly aligned we can only 50% of data, so ignored it.
- Problems are detected by EDA and addressed cleaning and validation process



Data Profile

Source: <https://data.cityofnewyork.us/Public-Safety/Motor-Vehicle-Collisions-Vehicles/bm4k-52h4>

Number of Tables: 3

Number of Records: 10,990,900 rows in total, number of columns ranges from 25 to 29

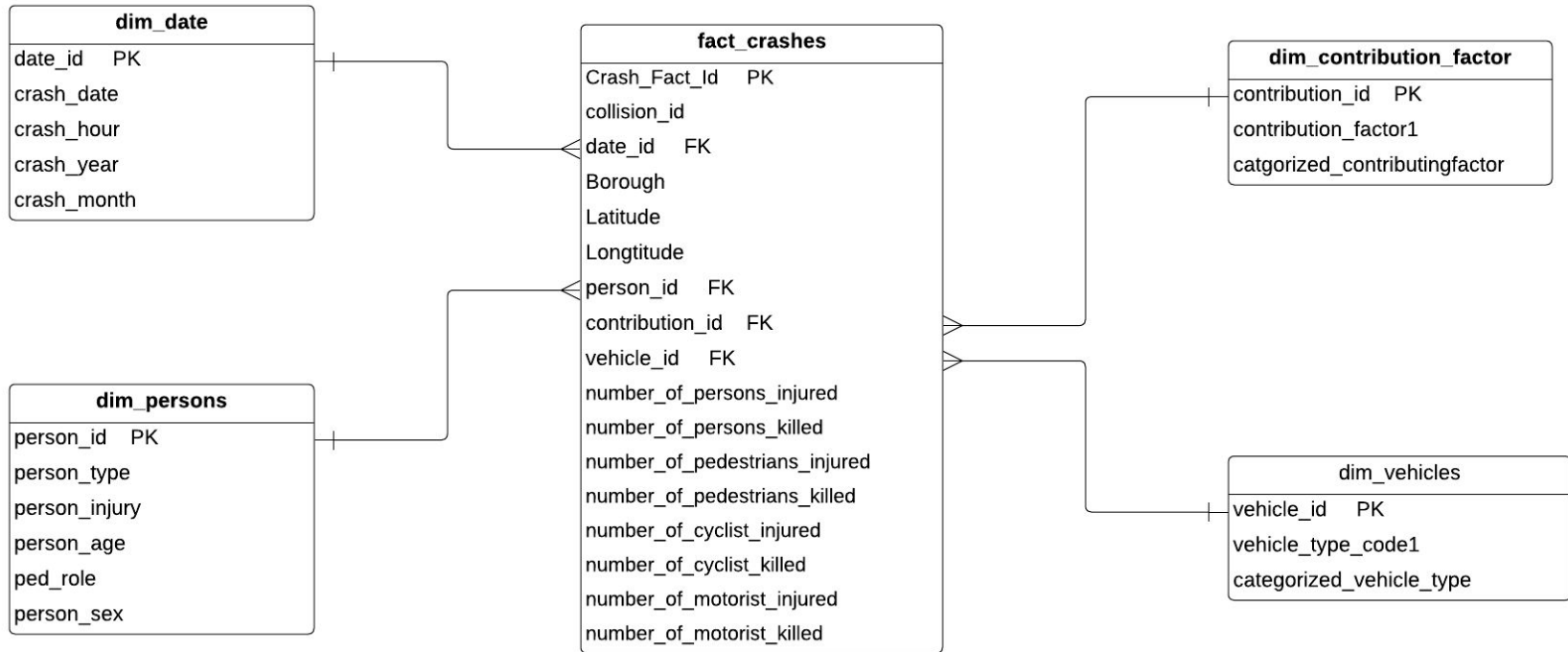
Data Type: Numeric, Text, Temporal, Geospatial

Data Acquisition Method: MySQL Connector through Python

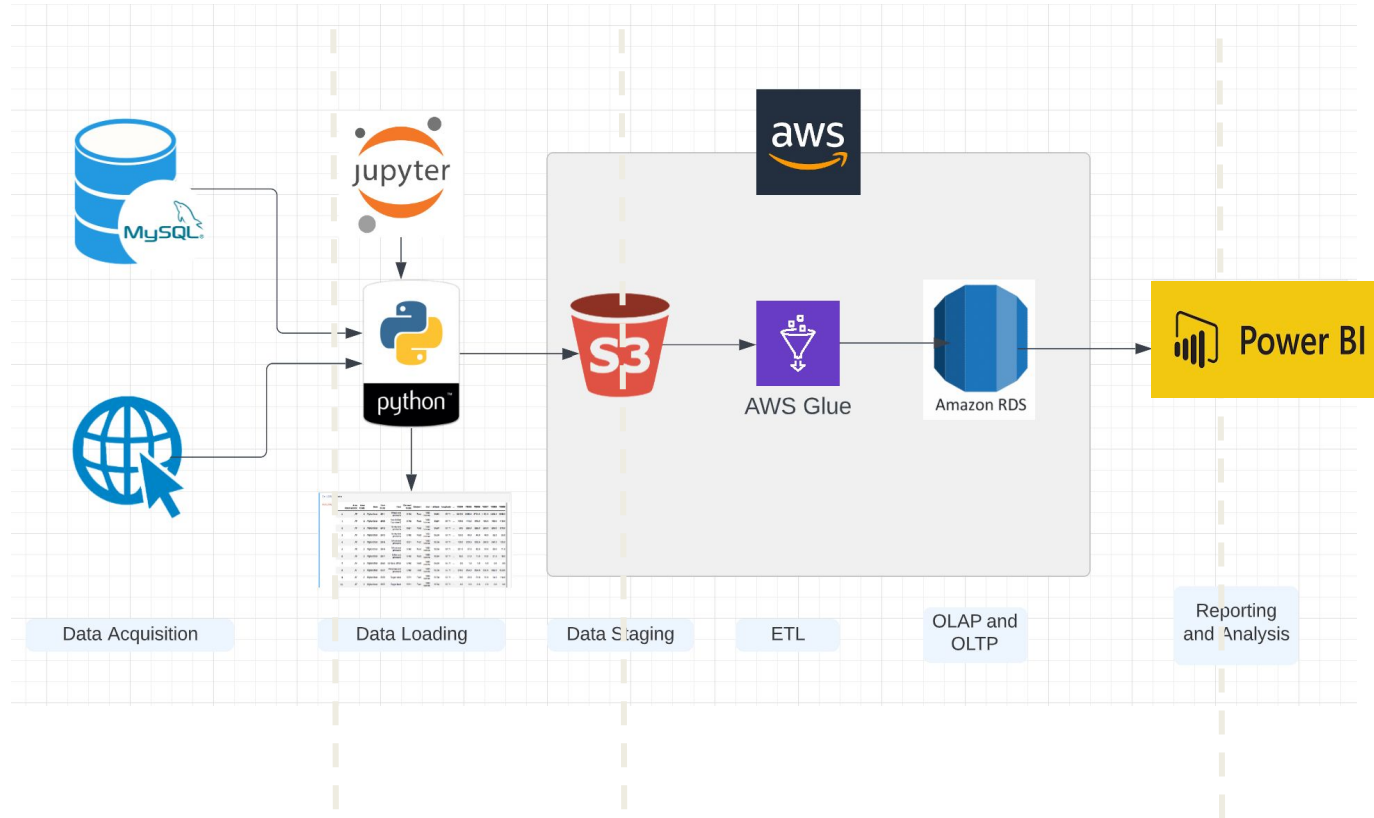
Bus Matrix

Facts	Date	Person	Contributing Factor	Vehicle	Location
Crash	X	X	X	X	X
Injury & Death	X	X		X	

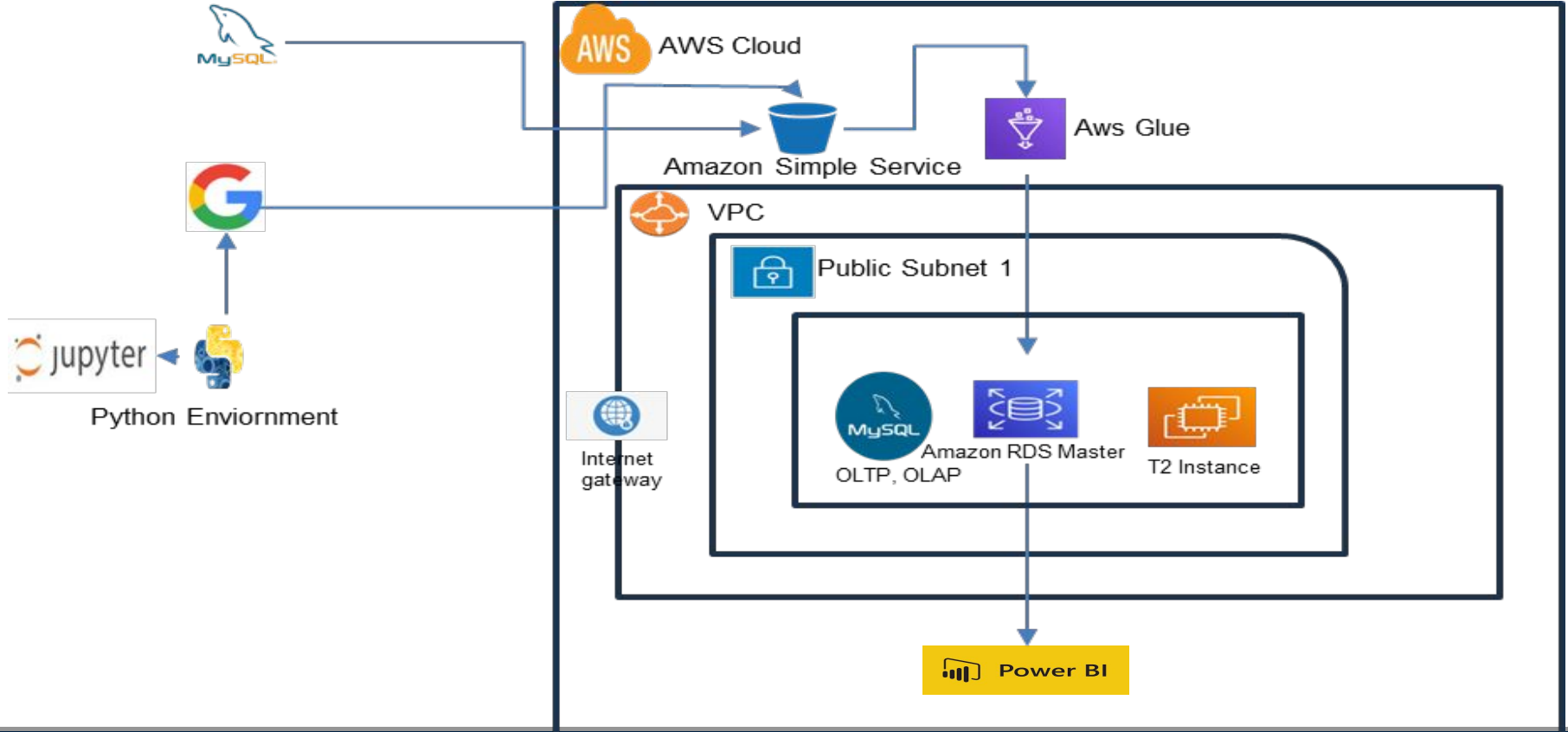
Dimensional Model



Work Flow



AWS Architecture



Dashboard

No of collisions

266.24K

People Killed

3856

People Injured

995K

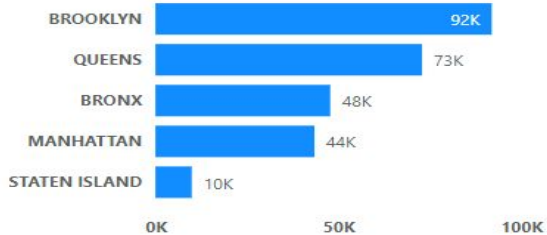
2020

2021

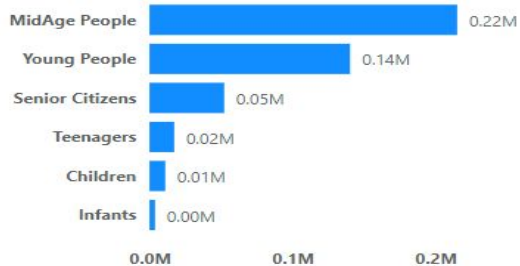
2022

2023

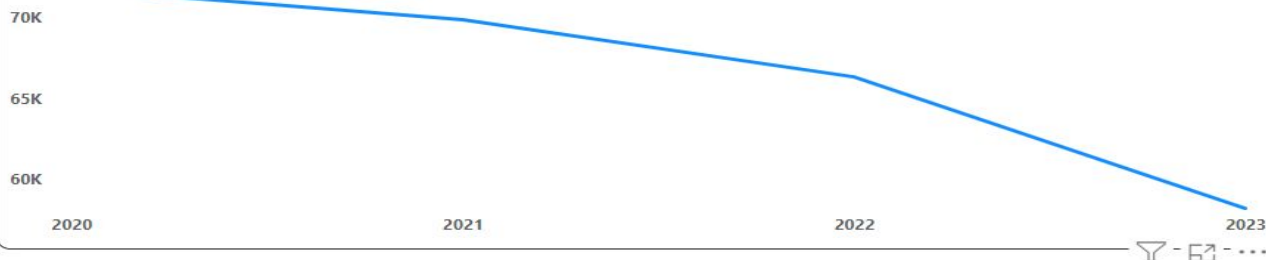
Borough wise Accidents



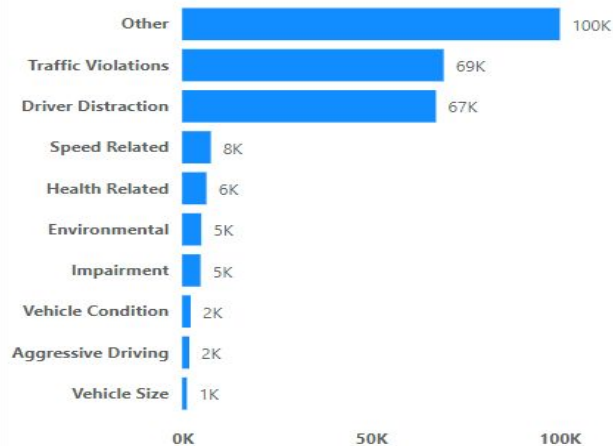
Age groups involved in Accidents



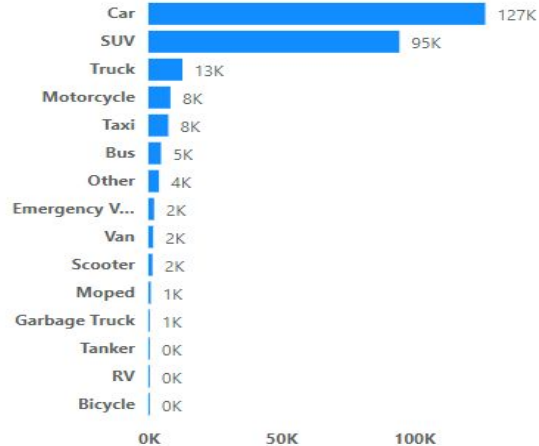
Year & month wise Accidents Trend



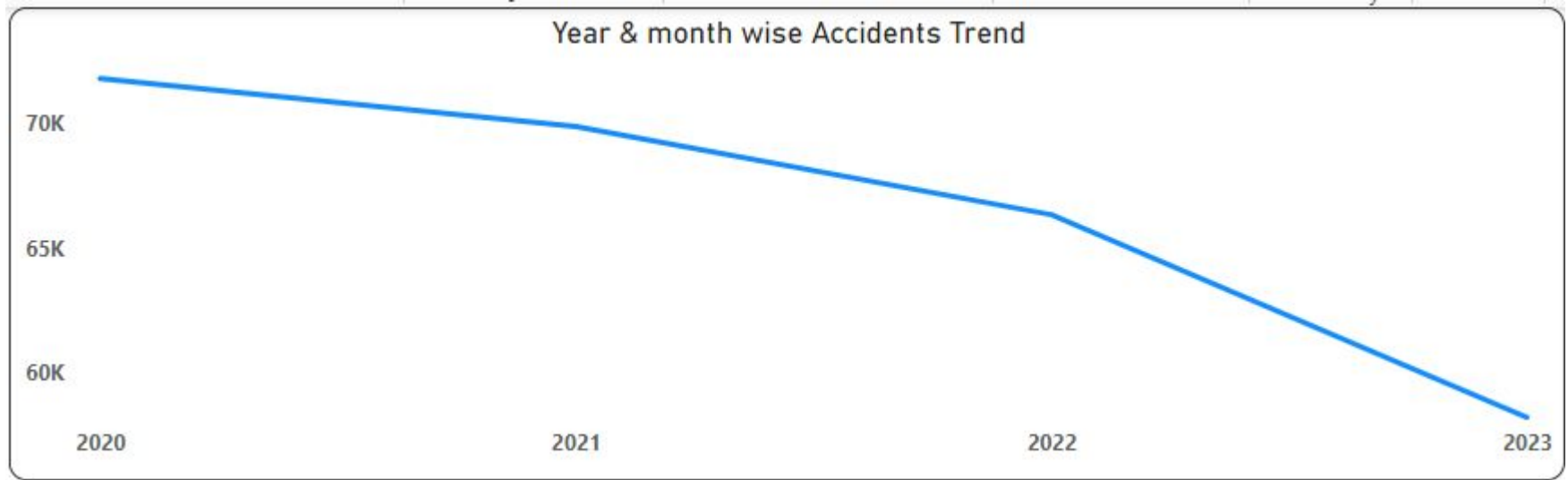
ContributingFactor for Accidents



Type of Cars involved in Accidents

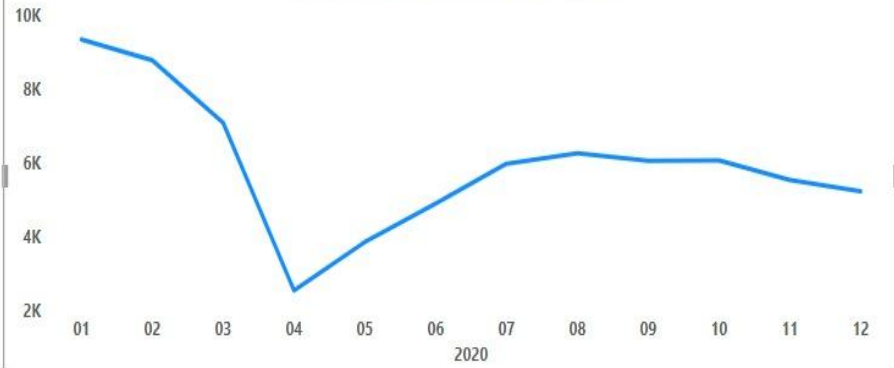


Visualization

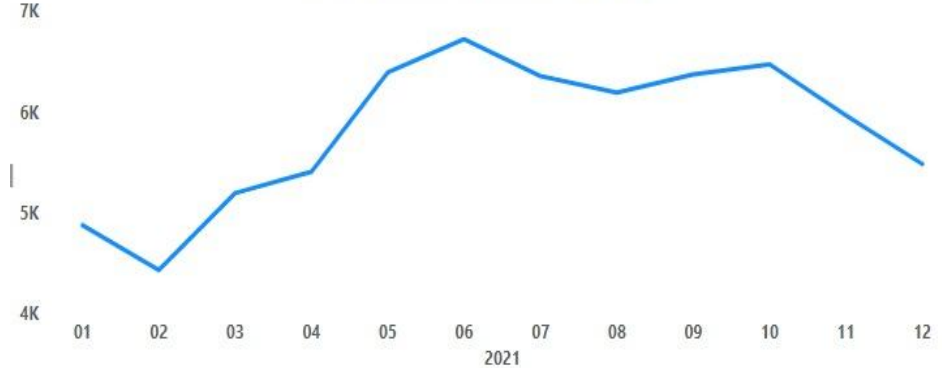


Visualization

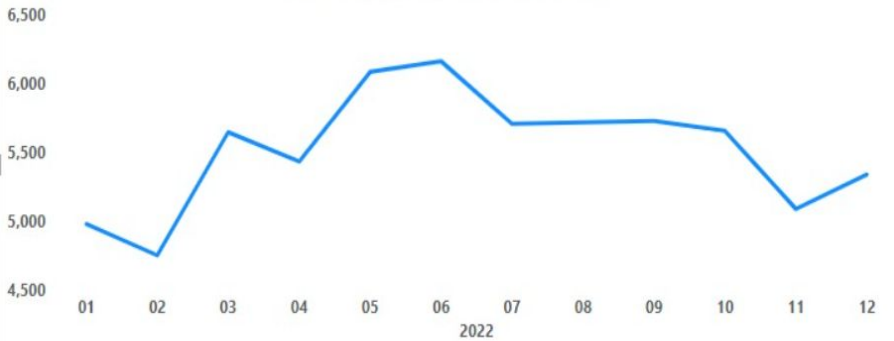
Year & month wise Accidents Trend



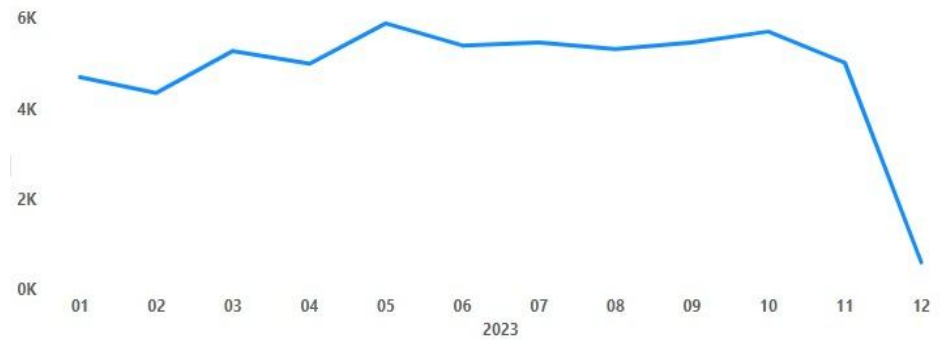
Year & month wise Accidents Trend

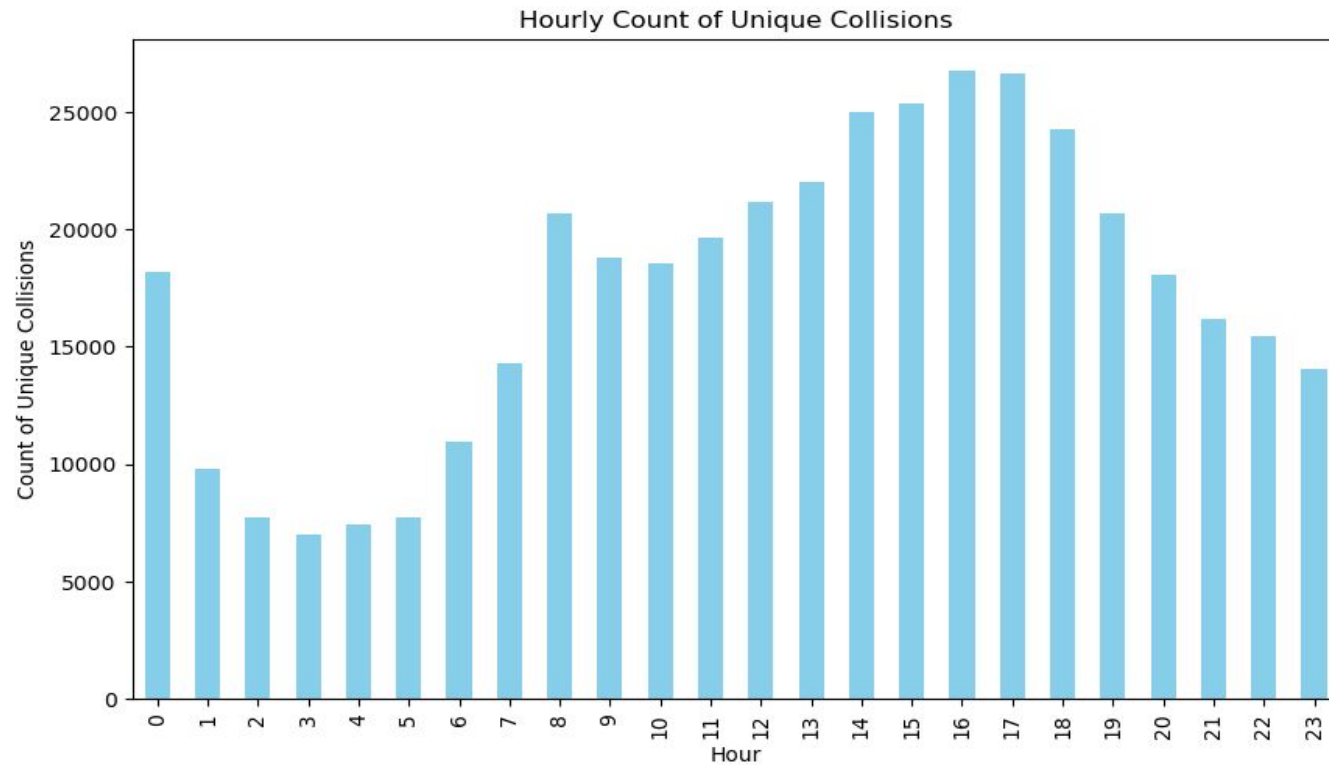


Year & month wise Accidents Trend



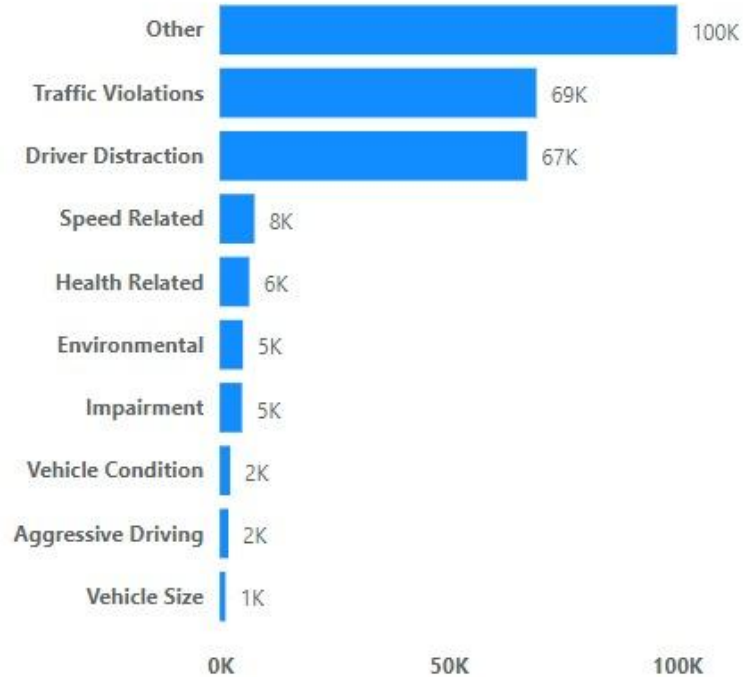
Year & month wise Accidents Trend



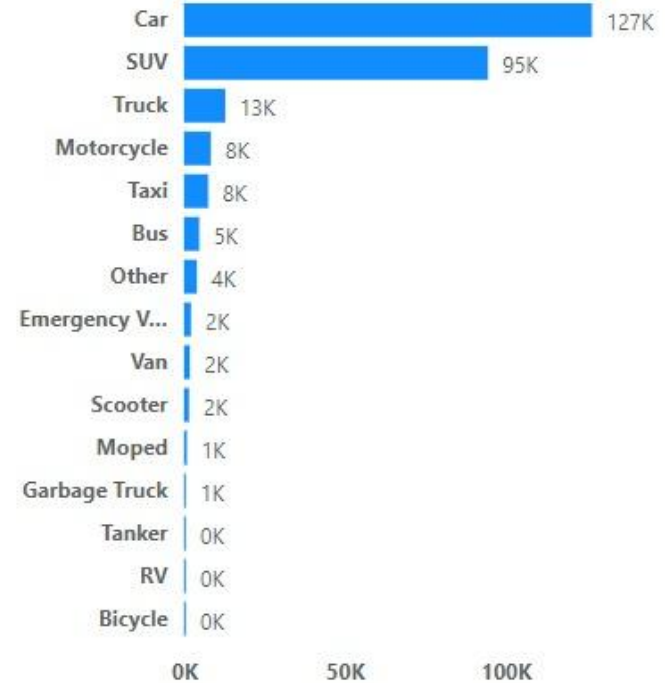


Visualization

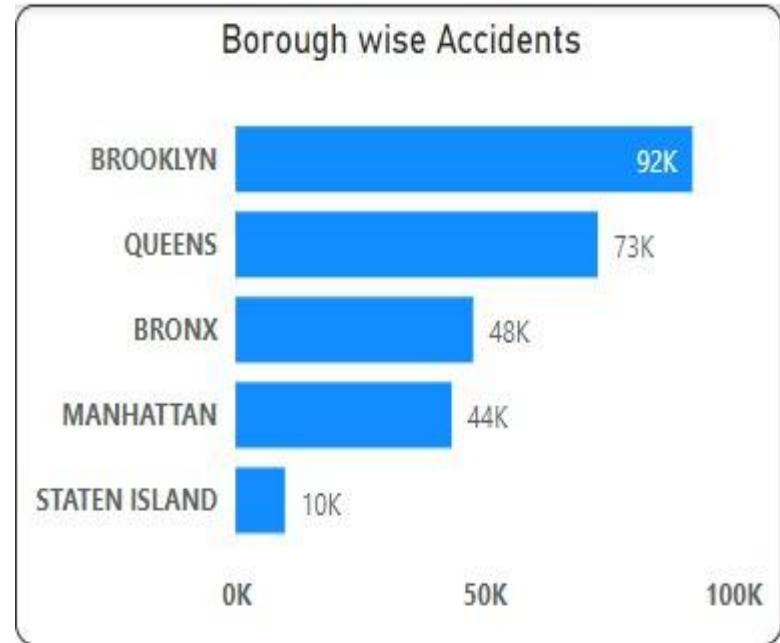
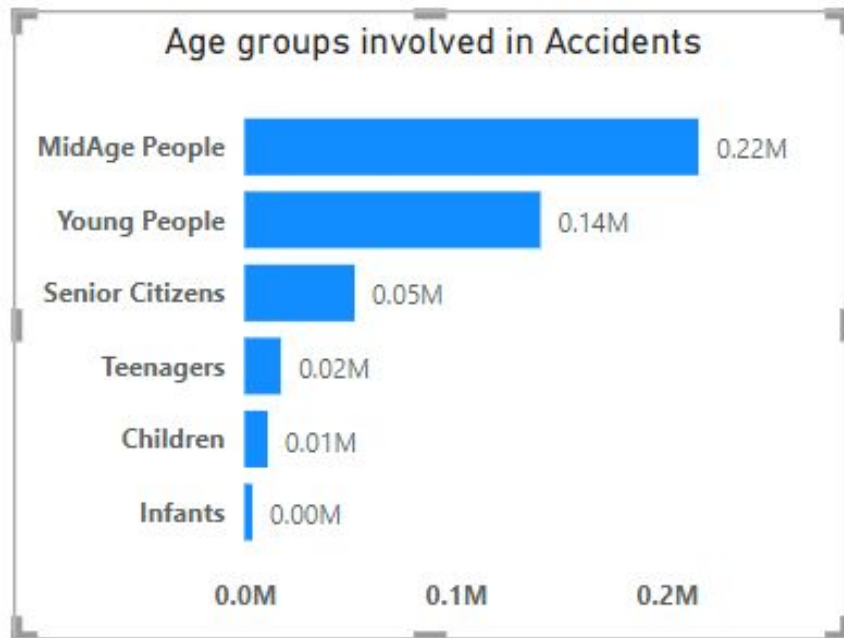
ContributingFactor for Accidents



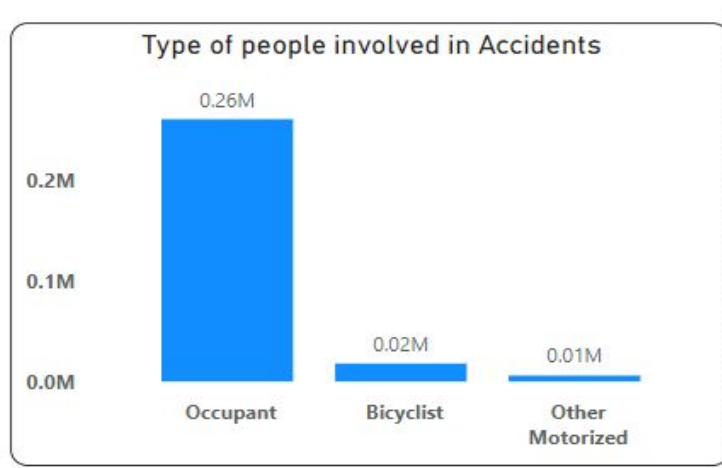
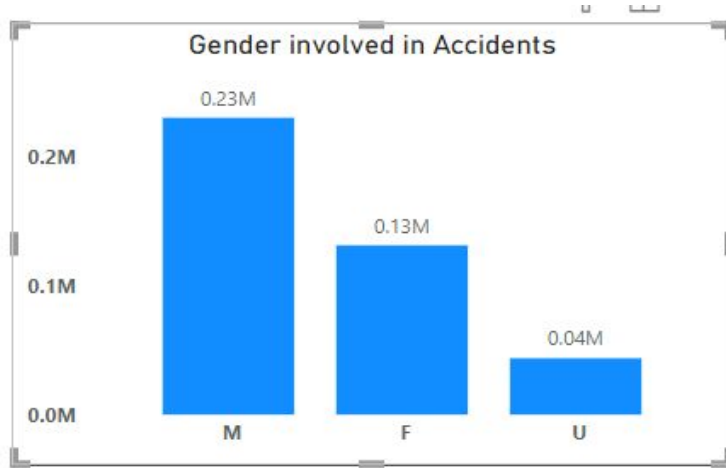
Type of Cars involved in Accidents



Visualization



Visualization



borough	Sum of Total Killed	Sum of Total Injured
BRONX	620	184988
BROOKLYN	1412	361836
MANHATTAN	589	128424
QUEENS	1019	278979
STATEN ISLAND	216	40680
Total	3856	994907

personsex	Sum of Total Killed	Sum of Total Injured
F	871	315635
M	2771	632828
U	214	46444
Total	3856	994907



Demo

S3 Buckets

Amazon S3 > Buckets > final-project-output-datasets-ia > Motor_Vehicle_Collisions_-_Crashes_20231209/

Motor_Vehicle_Collisions_-_Crashes_20231209/ [Copy S3 URI](#)

Objects | Properties

Objects (2) [Info](#)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

[Refresh](#) [Copy S3 URI](#) [Copy URL](#) [Download](#) [Open](#) [Delete](#) [Actions](#) [Create folder](#) [Upload](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	run-1703003330588-part-r-00000	-	December 19, 2023, 22:00:13 (UTC+05:30)	47.2 MB	Standard
<input type="checkbox"/>	run-1703003330588-part-r-00001	-	December 19, 2023, 22:00:13 (UTC+05:30)	23.8 MB	Standard

Motor_Vehicle_Collisions_-_Persons_20231209/ [Copy S3 URI](#)

Objects | Properties

Objects (5) [Info](#)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

[Refresh](#) [Copy S3 URI](#) [Copy URL](#) [Download](#) [Open](#) [Delete](#) [Actions](#) [Create folder](#) [Upload](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	run-1703003957495-part-r-00000	-	December 19, 2023, 22:09:45 (UTC+05:30)	42.0 MB	Standard
<input type="checkbox"/>	run-1703003957495-part-r-00001	-	December 19, 2023, 22:09:52 (UTC+05:30)	41.6 MB	Standard
<input type="checkbox"/>	run-1703003957495-part-r-00002	-	December 19, 2023, 22:09:51 (UTC+05:30)	41.5 MB	Standard
<input type="checkbox"/>	run-1703003957495-part-r-00003	-	December 19, 2023, 22:09:53 (UTC+05:30)	41.0 MB	Standard



Demo

AWS Glue

AWS Glue Studio [Info](#)

Create job [Info](#)

Author in a visual interface focused on data flow. [Visual ETL](#)

Author using an interactive code notebook. [Notebook](#)

Author code with a script editor. [Script editor](#)

► **Example jobs** [Info](#) [Create sample job](#)

Your jobs (3) [Info](#) [Refresh](#) [Actions](#) [Run job](#)

<input type="checkbox"/>	Job name	Type	Last modified	AWS Glue version
<input type="checkbox"/>	irs-990-job	Glue ETL	12/22/2023, 8:09:43 AM	4.0
<input type="checkbox"/>	Final_Project_glue_job	Glue ETL	12/20/2023, 5:09:56 AM	4.0
<input type="checkbox"/>	Final_Project_glue_job_Persons	Glue ETL	12/20/2023, 5:05:12 AM	4.0

Crawlers

Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

Crawlers (3) [Info](#) [Refresh](#) [Action](#) [Run](#) [Create crawler](#)

Last updated (UTC) December 22, 2023 at 20:27:00

View and manage all available crawlers.

<input type="checkbox"/>	Name	State	Schedule	Last run	Last run time...	Log	Table changes fr...
<input type="checkbox"/>	Crawler+persons	Ready		Succeeded	December 19, 20...	View log	1 created
<input type="checkbox"/>	irs-990-crawler	Ready		Succeeded	December 22, 20...	View log	1 created
<input type="checkbox"/>	s3-to-RDS	Ready		Succeeded	December 19, 20...	View log	-



Demo

[AWS Glue](#) > Tables

Tables

A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition.

Tables (3)

Last updated (UTC)
December 22, 2023 at 20:36:43



Delete

Add tables using crawler

Add table

View and manage all available tables.

Filter tables

< 1 >

<input type="checkbox"/>	Name ▲	Database ▼	Location ▼	Classification ▼	Deprecated ▼	View data	Data quality
<input type="checkbox"/>	irs_990_1	irs-990-2	s3://irs-990-1/	CSV	-	Table data	View data quality
<input type="checkbox"/>	motor_vehicle_collisio	crashes_databaes	s3://final-project-outp	CSV	-	Table data	View data quality
<input type="checkbox"/>	motor_vehicle_collisio	crashes_databaes	s3://final-project-outp	CSV	-	Table data	View data quality



Assumptions

Data Completeness and Accuracy: We assumed that the crash data captured in the fact table, along with the associated dimension tables (dim_date, dim_persons, dim_vehicles, dim_contribution_factor), is complete and accurate.

Stable and Consistent Data Sources: The project assumes that the data sources providing the information for the crash data are stable and consistent over time.

Policy Impact and Relevance Assumption: We assumed that the data and insights obtained from the fact table and related dimension tables will be directly relevant and actionable for policy makers and business stakeholders.

Challenges

- i) Finding data
- ii) Solving error from AWS connection

Thank You for Listening!

