

# NYC Taxi and Citi Bike Data Analysis

## Project Report

**Data - 228 (Big Data Technologies and Applications)**

Team Dataluation

Abinaya Seshadre, Neelima Jagtap, Rishikumar Ravichandran, Sruthi Mallarapu

### **I. Abstract**

Yellow taxis have been a well-renowned symbol and phenomenon of New York city. They serve the 5 major boroughs of NYC: Manhattan, Queens, Bronx, Brooklyn, and Staten Island. NYC TLC is the government agency that oversees the operations of taxi services in NYC including controlling and handing out licenses to all taxi businesses, even to competitors to yellow taxis like Uber. Before the pandemic hit, there were about 1 million taxi trips made every day in NYC. TLC collects the trips data, stores them in its repository and shares those for public use.

Citi bike is the bike sharing program of NYC. These bikes are available throughout the day, throughout the year and can be picked up from any station and returned any of the docking locations, which is very convenient for commuters with various purposes and generate tons of transaction data. Currently, there are about 21500 bikes with 1400+ stations in and around the major boroughs.

### **II. Project Goals**

We had the following goals in mind for this term project:

- Find and use big data - datasets with million rows of data.
- Implement data modelling and data processing for data analysis and visualizations.
- Explore and use cloud technology for all our requirements.
- Uncover exciting insights with analysis and visualizations.

Later sections will describe, in detail, how each of our project goals were achieved.

### **III. Motivation**

NYC TLC vehicles make about 100000 trips every day. This generates huge data for TLC that can be analyzed to answer a multitude of questions.

1. Identify hotspots.
2. Identify various trends.
3. Compare taxi and bike data.
4. Capture the effects of the pandemic.

## **IV. Methodology**

This section gives a detailed description of the methods we followed to accomplish our project goals.

The first step we took was that of requirements gathering which consisted of the following steps:

### **Project Planning**

- Probing the internet for open datasets.
- Discussing about our goals and requirement of the project.
- Fixing the datasets
- Choosing our tools
- Choosing our methodology
- Following class lectures, homework, and online resources to understand concepts and help with the project implementation.

We ended up with the following high-level decisions, after which we came up with our end-to-end data flow chart.

- We chose to do our project using AWS services.
- We decided our datasets and took up combined analysis of Citi bike and Taxi data in NYC.
- We also wanted to showcase our NYC TLC Data Analysis for insights on peak hour, areas, and times of services and COVID effects.

The below flow chart depicts the process we followed for data ingestion from the data sources to our pipelines and to the final stages of data analysis and visualization.

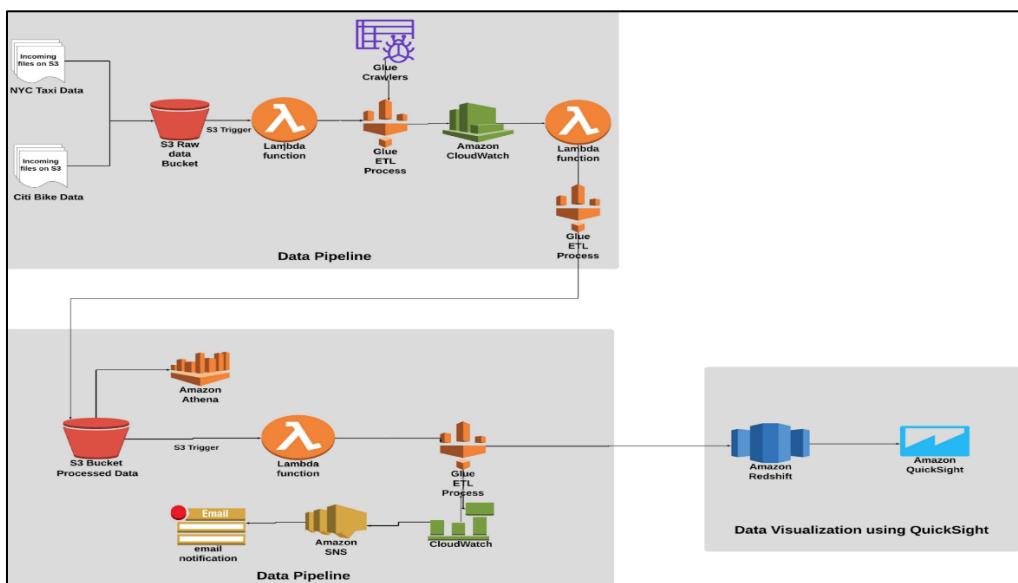


Figure 1: Data flow from source to destination via pipelines

## Process Overview

The below steps give an overview of the processes undertaken for data ingestion to our pipelines:

- Imported Data to S3 buckets using **AWS CLI**.
- Created an **S3 Trigger** that invoked the **Lambda function** which crawls the data in S3.
- Created **Cloud Watch Event** that triggers another Lambda function which perform the **ETL job** and stores the processed data in S3.
- Created an S3 Trigger in processed bucket that invoked the Lambda function that performed the ETL job of transferring processed data from S3 to **Redshift**.
- Created Cloud Watch Event which invokes **AWS SNS** and sends email notification on successful execution.
- Connected Redshift to QuickSight and created visualization dashboards.

## Process Description

### **Data sources**

Our data source for NYC Taxi was from the public repository available in public S3 buckets where each month's data is loaded in batches after the month is over. Instead of taking the same years data, we chose data that spanned the months of March to May for the years 2019 and 2020 so we will be able to capture the effects of the pandemic on our chosen business. The chosen datasets had 26.6 million data rows in total.

The datasets had details about the taxi trips including fields like pickup location, drop location, timestamp, fares, distance and so on.

Our data for New York Citi bikes was procured from NYC Open data S3 buckets. This dataset contains data about the pickup location, drop location, timestamp, duration, distance and so on. For Citi bike, we chose data that spanned the months of March to May for the year 2019, for comparison with our taxi data. The chosen dataset had 5 million data rows in total.

We imported the NYC Yellow taxi data and Citi bike data to our project S3 buckets using AWS CLI.

### **Data Storage**

Amazon Simple Storage Service, S3 acts as a file storage system and provides storage space for various data and objects. S3 has scalable infrastructure and is a public resource in the cloud provided by Amazon Web Services. The advantages of using S3 bucket is storing all data in one location with access to any of the services and tools provided by AWS. Also, the scalability factor makes it more beneficial for the implementation of projects.

Our data files were loaded into our S3 bucket for easy storage and further loading into Amazon Redshift. We decided on Amazon S3 due to its better compatibility with Amazon Redshift and QuickSight.

### **Cleansing (Data Pre-processing)**

Raw data almost always has redundancies and gaps. Such inconsistencies should be removed by either deletion or replacement before further processing can be done so the results generated from the data are accurate and not misleading. Cleansing can be accomplished by using Glue scripts. We did not want to go for other online data Cleansing tools as we wanted to experiment and achieve as much as possible with the cloud solutions and AWS.

Amazon Glue is a serverless processing and computing service provided by AWS. It is useful for data integration and formation of metadata for table data. It acts as a crawler which crawls through the data to detect the schema and composes the metadata for the data, after which data is loaded into the chosen target. Data can be pre-processed using Glue. Glue provides with a data catalog, which could be used to analyze and understand our data better. It is an Extract Transform Load (ETL) and event driven service to build our data in a better shape. The data is ingested into Glue with the help of a crawler and a job.

### **Data Loading**

All the data, once cleansed, must be loaded into a common repository for easy access. We had our raw data and processed data in S3 and loaded our processed data to Redshift for further analysis using SQL queries.

### **Data Modeling**

After the necessary, processed data are available, we used tools like draw.io to identify the entities, attributes, relationships and generate the data model, as the AWS data modeler available as part of DynamoDB was not available in the free tier. This data model is important to understand the fields in the data sets and relationships between the entities, if available and would also assist us in generating novel ideas of performing analysis on the data.

### **Data Analysis**

Various types of analytics, including ‘Descriptive analysis, Exploratory analysis and Predictive analysis’ can be carried out on the data loaded into our AWS database, Amazon Redshift.

Amazon Redshift is a fully managed, cloud-based data warehouse provided by AWS. It is SQL compatible and SQL queries are used to perform DDL and DML operations on the tables and the large datasets. Redshift can store and operate (process joins) on even petabyte sized data.

We performed data analysis using SQL queries in Redshift and later also in QuickSight using graphs and charts.

### **Visualization**

Finally, the results were presented in a graphical manner for easy and comprehensive visualization by using Amazon’s visualization tool, QuickSight.

Amazon QuickSight is AWS's visualization tool. It is highly scalable and consists of business intelligence and machine learning capabilities built in to easily create insights, visuals

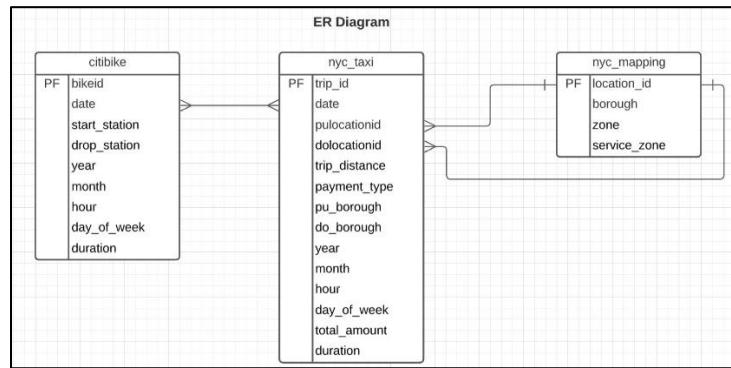
and dashboards for our applications, data models and projects.

QuickSight provides access to load data from many sources like S3, Redshift which are AWS services and sources external to AWS like PostgreSQL, Oracle, Teradata and so on. We can create graphs and charts using our own custom SQL queries.

## V. Implementation

This section dives deeper into our implementation steps and elaborates more on our described methodologies. It includes a detailed description of every step with appropriate screenshots showing execution in AWS.

### Our Data Model



Our data model has the entities for the Yellow Taxi and Citi Bikes data tables and also describes the attributes required in each of the entity.

### Implementation Steps:

#### **Data Ingestion**

#### **Automated ETL pipeline**

We automated the ETL pipeline from creating the data catalog to transferring the processed data to Redshift. Whenever new data is uploaded into our S3 buckets, below steps are automatically triggered:

1. Creation of the Data Catalog
2. ETL job Execution
3. Data transfer to Redshift.
4. Email Notification upon Successful Transfer.

#### **1. Creating the Data Catalog**

To determine the schema of the raw data in S3, data is crawled using AWS Glue. This creates a data catalog which has the schema of the data. We automated this step by

creating a trigger in S3 bucket that invokes the AWS Lambda function. This function starts the AWS Glue crawler and catalogs the data.

## 2. ETL JOB Execution

Upon Successful completion of Crawler, the second AWS Lambda function is triggered by using Amazon CloudWatch Events rule. This function invokes the ETL job.

The ETL job is created using AWS Glue studio that cleans the data and transfers it to processed S3 bucket. Cleaning process involved dropping unnecessary columns, changing data types, filtering the data, and performing the joins with mapping files.

## 3. Transferring the data to Redshift.

A Trigger is created in the S3 bucket with processed data that invokes the third AWS Lambda function. This function triggers another ETL job that transfers the data from S3 to Redshift. This ETL job is created using AWS Glue and as a pre-requisite, we also configured the Redshift connection and created necessary, predefined tables in Redshift.

## 4. Email Notification upon Successful Transfer.

Created a rule in CloudWatch Events rule in such a way that when the status of the AWS Glue ETL job is ‘Success’, it sends an email notification using an SNS topic. Configured and set the email field for AWS to know which email IDs, notifications need to be sent in SNS Topic.

### Step by Step with screenshots

#### Data Ingestion into S3 buckets.

We first created separate buckets in S3 to store our raw Yellow taxi data and Citi Bike data and then loaded the data files using the CLI.

Objects (6)						
Objects are the fundamental entities stored in Amazon S3. You can use <a href="#">Amazon S3 inventory</a> to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. <a href="#">Learn more</a>						
<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class	Actions
<input type="checkbox"/>	<a href="#">yellow_tripdata_2019-03.csv</a>	csv	May 15, 2021, 15:29:13 (UTC-07:00)	692.6 MB	Standard	<input type="button" value="Copy URL"/> <input type="button" value="Open"/> <input type="button" value="Download"/> <input type="button" value="Delete"/> <input type="button" value="Actions"/> <input type="button" value="Create folder"/> <input type="button" value="Upload"/>
<input type="checkbox"/>	<a href="#">yellow_tripdata_2019-04.csv</a>	csv	May 15, 2021, 15:20:21 (UTC-07:00)	657.3 MB	Standard	<input type="button" value="Copy URL"/> <input type="button" value="Open"/> <input type="button" value="Download"/> <input type="button" value="Delete"/> <input type="button" value="Actions"/> <input type="button" value="Create folder"/> <input type="button" value="Upload"/>
<input type="checkbox"/>	<a href="#">yellow_tripdata_2019-05.csv</a>	csv	May 15, 2021, 15:20:21 (UTC-07:00)	669.0 MB	Standard	<input type="button" value="Copy URL"/> <input type="button" value="Open"/> <input type="button" value="Download"/> <input type="button" value="Delete"/> <input type="button" value="Actions"/> <input type="button" value="Create folder"/> <input type="button" value="Upload"/>
<input type="checkbox"/>	<a href="#">yellow_tripdata_2020-03.csv</a>	csv	May 15, 2021, 15:20:21 (UTC-07:00)	265.4 MB	Standard	<input type="button" value="Copy URL"/> <input type="button" value="Open"/> <input type="button" value="Download"/> <input type="button" value="Delete"/> <input type="button" value="Actions"/> <input type="button" value="Create folder"/> <input type="button" value="Upload"/>
<input type="checkbox"/>	<a href="#">yellow_tripdata_2020-04.csv</a>	csv	May 15, 2021, 15:20:21 (UTC-07:00)	20.7 MB	Standard	<input type="button" value="Copy URL"/> <input type="button" value="Open"/> <input type="button" value="Download"/> <input type="button" value="Delete"/> <input type="button" value="Actions"/> <input type="button" value="Create folder"/> <input type="button" value="Upload"/>
<input type="checkbox"/>	<a href="#">yellow_tripdata_2020-05.csv</a>	csv	May 15, 2021, 15:20:21 (UTC-07:00)	30.2 MB	Standard	<input type="button" value="Copy URL"/> <input type="button" value="Open"/> <input type="button" value="Download"/> <input type="button" value="Delete"/> <input type="button" value="Actions"/> <input type="button" value="Create folder"/> <input type="button" value="Upload"/>

Figure 2: Yellow Taxi Data of 6 months period (Mar2019-May2019, Mar2020-May2020) uploaded to our S3 bucket - **yellow-taxi-data-ny**.

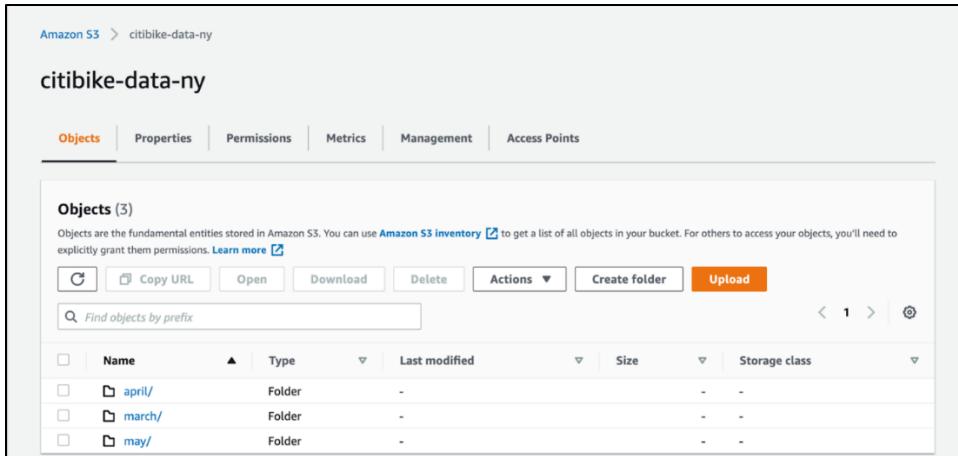


Figure 3: Citi bike data of three months period (Mar2019-May2019) has been uploaded into another S3 bucket - **citibike-data-ny**.

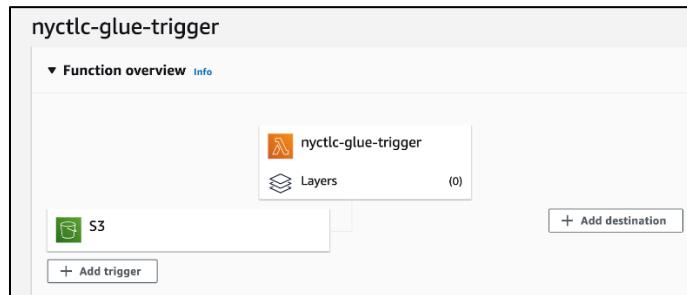
## Data Pipelines

We created the demo pipeline which we used for our project presentation and replicated this pipeline on our raw data buckets. Attached are the screenshots for the step-by-step process for both our project as well as demo pipeline.

Created the S3 trigger name **nyctlc-crawler-trigger** in the “**yellow-taxi-data-ny**” bucket. This trigger invokes the lambda function named “**nyctlc-glue-trigger**” whenever new file is uploaded in the bucket.



The Lambda function invoked above is used to create the crawler named “**datacrawler\_raw**” and then creates the data catalog.



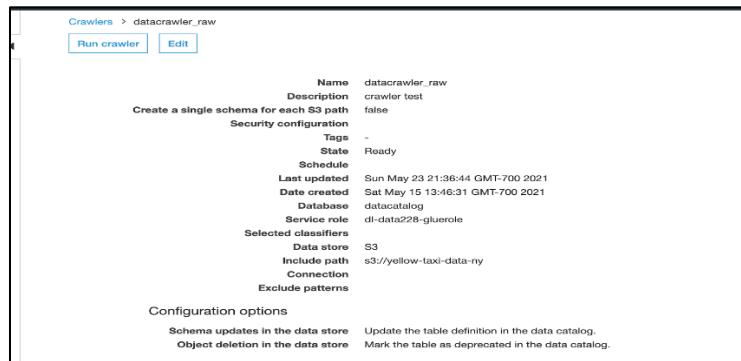
Code Snapshot for the above lambda function

```

var AWS = require('aws-sdk');
var glue = new AWS.Glue();
exports.handler = function(event, context,callback) {
    console.log(JSON.stringify(event, null, 3));
    var dbName = 'datacatalog';
    var params = [
        DatabaseInput: [
            Name: dbName,
            Description: 'project',
        ]
    ];
    glue.createDatabase(params, function(err, data) {
        var params1 = {
            DatabaseName: dbName,
            Name: 'datacrawler_raw',
            Role: 'service-role/dl-data228-gluerole',
            Targets: {
                S3Targets: [{ Path: 's3://yellow-taxi-data-ny' }]
            },
            Description: 'crawler test'
        };
        glue.createCrawler(params1, function(err1, data1) {
            startCrawler('datacrawler_raw', function(err2,data2){
                if(err2) callback(err2)
                else callback(null,data2)
            })
        });
    });
    function startCrawler(name,callback){
        var params = {
            Name: name,
        };
        glue.startCrawler(params, function(err, data) {
            callback(null, data)
        });
    }
});

```

Next, **datacrawler\_raw** crawler is created by the invocation of the above lambda function.



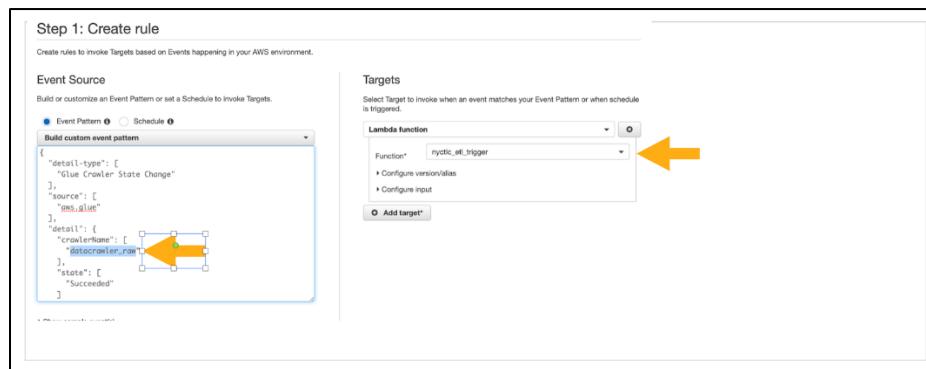
After the datacrawler\_raw crawls the data, "yellow\_taxi\_data\_ny" catalog table is created.

Name	yellow_taxi_data_ny
Description	demo
Database	demo
Classification	csv
Location	s3://yellow-taxi-data-ny/
Connection	
Deprecated	No
Last updated	Sun May 23 21:33:54 GMT-700 2021
Input format	org.apache.hadoop.mapred.TextInputFormat
Output format	org.apache.hadoop.hive.qi.io.HiveIgnoreKeyTextOutputFormat
Serde serialization lib	org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
Serde parameters	field.delim : , skip.header.line.count : 1 sizeKey : 2448540037 objectCount : 6
Table properties	UPDATED_BY_CRAWLER : nyctic-rawdata-crawler CrawlerSchemaDeserializerVersion : 1.0 recordCount : 17187518 averageRecordSize : 142 CrawlerSchemaDeserializerVersion : 1.0 compressionType : none columnsOrdered : true areColumnsQuoted : false delimiter : , typeOfData : file
Schema	

Schema of the raw data

Column name	Data type	Partition key
1 vendorid	bigint	
2 tpep_pickup_datetime	string	
3 tpep_dropoff_datetime	string	
4 passenger_count	bigint	
5 trip_distance	double	
6 ratecodeid	bigint	
7 store_and_fwd_flag	string	
8 pulocationid	bigint	
9 dolocationid	bigint	
10 payment_type	bigint	
11 fare_amount	double	
12 extra	double	
13 mta_txs	double	
14 tip_amount	double	
15 tolls_amount	double	
16 improvement_surcharge	double	
17 total_amount	double	
18 congestion_surcharge	double	

An event rule is created in Cloud Watch that invokes another lambda function upon successful crawling of the data.



Code snapshot for the second lambda function “**nyc\_tlc\_etl\_trigger**”. It invokes the ETL job named ‘**Automated\_Job**’.

```

 1 var AWS = require('aws-sdk');
 2 var sns = new AWS.SNS({ region: "us-east-1" });
 3 var s3 = new AWS.S3();
 4 var glue = new AWS.Glue({apiVersion: '2017-03-31'});
 5 exports.handler = function(event, context, callback) {
 6     var params = {
 7         JobName: 'Automated_Job',
 8         Timeout: 20,
 9     };
10     glue.startJobRun(params, function(err1, data1) {
11         if (err1) {
12             console.log(err1, err1.stack);
13         } else {
14             console.log(data1);
15         };
16         console.log(JSON.stringify(event, null, 3));
17     });
18 }

```

The screenshot shows the AWS Lambda function editor for the 'nyc\_tlc\_etl\_trigger' function. The code is written in JavaScript and uses the AWS SDK to interact with AWS services like SNS, S3, and Glue. It triggers an AWS Glue job named 'Automated\_Job' with a timeout of 20 seconds. An orange arrow points from the 'JobName' parameter in the code towards the 'Automated\_Job' entry in the AWS Lambda function configuration.

The ETL job is created using AWS Glue studio that cleans the data and transfers it to processed S3 bucket. Cleaning process involved dropping unnecessary columns, changing data types, filtering the data, and performing the joins with mapping files. Mapping files are present in S3 bucket named “**mapping-file-nyctlc-nvirginia**”.

Amazon S3 > mapping-file-nyctlc-nvirginia

### mapping-file-nyctlc-nvirginia

Publicly accessible

**Objects (1)**

mapping.csv

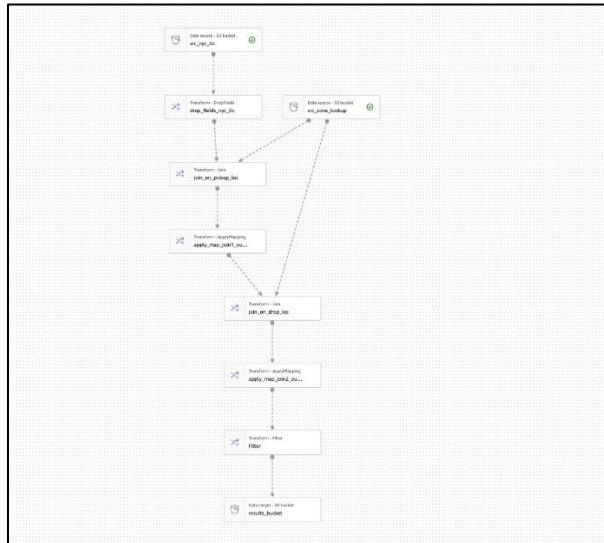
Last modified: May 15, 2021, 10:55:49 (UTC-07:00)

Type: csv

Size: 10.5 KB

Storage class: Standard

## ETL FLOW



Processed Data is stored in S3 processed buckets.

Source Transform Target Undo Redo Remove

Node properties Data target properties - S3 Output schema

Format: CSV

Compression Type: None

S3 Target Location: s3/yellow-taxi-data-ny-processed/

Data Catalog update options:

- Do not update the Data Catalog
- Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions
- Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions

Partition keys - optional

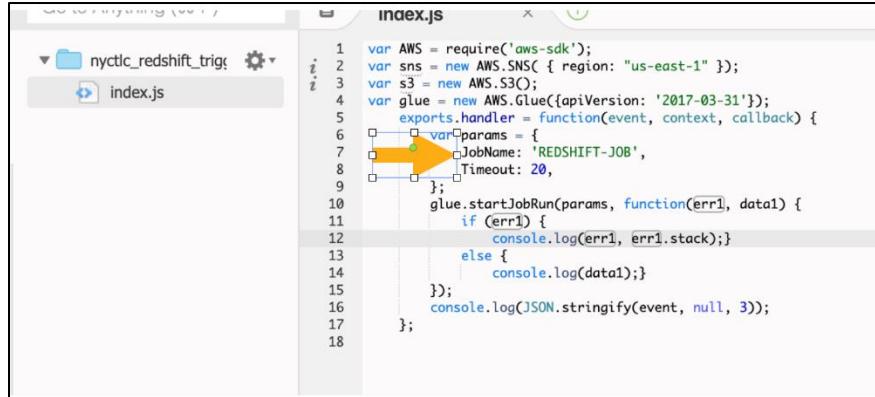
The cleaned data is partitioned and is stored in a processed S3 bucket named “**yellow-taxi-data-ny-processed**”.

Objects (120)						
Objects are the fundamental entities stored in Amazon S3. You can use <a href="#">Amazon S3 Inventory</a> to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. <a href="#">Learn more</a>						
		Type	Last modified	Size	Storage class	
<input type="checkbox"/>	run-1621118310587-part-r-00026	-	May 15, 2021, 15:51:47 (UTC-07:00)	786.2 KB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00029	-	May 15, 2021, 15:54:52 (UTC-07:00)	63.4 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00030	-	May 15, 2021, 15:55:25 (UTC-07:00)	85.2 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00031	-	May 15, 2021, 15:55:03 (UTC-07:00)	121.6 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00032	-	May 15, 2021, 15:55:11 (UTC-07:00)	71.8 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00033	-	May 15, 2021, 15:53:15 (UTC-07:00)	38.6 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00034	-	May 15, 2021, 15:53:03 (UTC-07:00)	25.5 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00035	-	May 15, 2021, 15:52:04 (UTC-07:00)	6.8 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00036	-	May 15, 2021, 15:54:00 (UTC-07:00)	71.5 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00037	-	May 15, 2021, 15:53:30 (UTC-07:00)	35.9 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00038	-	May 15, 2021, 15:51:52 (UTC-07:00)	1.9 MB	Standard	
<input type="checkbox"/>	run-1621118310587-part-r-00039	-	May 15, 2021, 15:53:44 (UTC-07:00)	50.9 MB	Standard	

We have used Athena service to check if the processed data is in correct format. Created the S3 trigger name **trigger-Redshiftjob** in the “yellow-taxi-data-ny-processed” bucket. This trigger invokes the “**nyctlc\_redshift\_trigger**” lambda function as soon as the processed data comes.

Event notifications (1)				
Send a notification when specific events occur in your bucket. <a href="#">Learn more</a>				
Name	Event types	Filters	Destination type	Destination
trigger-Redshiftjob	All object create events	-	Lambda function	nyctlc_redshift_trigger

This Lambda function invokes “**REDSHIFT-JOB**”.



```

var AWS = require('aws-sdk');
var sns = new AWS.SNS({ region: "us-east-1" });
var s3 = new AWS.S3();
var glue = new AWS.Glue({apiVersion: '2017-03-31'});
exports.handler = function(event, context, callback) {
  var params = {
    JobName: 'REDSHIFT-JOB',
    Timeout: 20,
  };
  glue.startJobRun(params, function(err1, data1) {
    if (err1) {
      console.log(err1, err1.stack);
    } else {
      console.log(data1);
    };
    console.log(JSON.stringify(event, null, 3));
  });
}

```

Configured the redshift Connection that is required by the **REDSHIFT-JOB** to transfer the data from processed S3 bucket to Redshift.

The screenshot shows the 'Connections' section in the AWS Management Console. A connection named 'Redshift\_connection' is selected. The details shown include:

- Type: JDBC
- JDBC URL: jdbc:redshift://redshift-cluster-1.cqk8mnrrhted.us-east-1.redshift.amazonaws.com:5439/dev
- VPC Id: vpc-5d58d020
- Subnet: subnet-2d80b623
- Security groups: sg-baa054be
- Require SSL connection: false
- Description: -
- Username: awsuser
- Created: 15 May 2021 4:17 PM UTC-7
- Last modified: 15 May 2021 4:17 PM UTC-7

Successful execution of **REDSHIFT-JOB**. It means data has transferred to Redshift tables.

The screenshot shows the 'Jobs' section in the AWS Glue console. The 'REDSHIFT-JOB' is listed as successful. The execution details are as follows:

Run ID	Retry attempt	Run status	Error	Output Logs	Error logs	Glue version	Maximum capacity	Triggered by	Start time	End time	Start-up time	Execution time	Timeout	Delay	Job run input
jr_916052e2...	-	Succeeded		Logs	Error logs	2.0	10		15 ...	15 ...	18 secs	5 mins	2880 mins		s3://aws-glu...

Another event rule is created in Cloud Watch that invokes SNS services. This service notifies through an email upon successful completion of the **REDSHIFT-JOB**.

The screenshot shows the 'Event Rules' section in the AWS CloudWatch Events console. A rule is configured to trigger on the 'REDSHIFT-JOB' success event. The target is an SNS topic named 'demostack-SNSProcessedEvent-10QL7FW4JPSV'. The rule details are as follows:

**Summary**

ARN: arn:aws:events:us-east-1:577417548496:rule/demostack-Opseventrule-1B9RLHBNUEMRM

Event pattern:

```
{
  "detail-type": [
    "Glue Job State Change"
  ],
  "source": [
    "aws:glue"
  ],
  "detail": [
    "jobName": [
      "REDSHIFT-JOB"
    ],
    "state": [
      "SUCCEEDED"
    ]
  ]
}
```

Status: Enabled

Description: EventRule

Monitoring: Show metrics for this rule

**Targets**

Type	Name	Input	Role	Additional parameters
SNS topic	demostack-SNSProcessedEvent-10QL7FW4JPSV	Matched event		

Configured Email in AWS SNS.

Details									
Name <b>demostack-SNSProcessedEvent-10QL7FW4JP5V</b>	Display name -								
ARN arn:aws:sns:us-east-1:577417548456:demostack-SNSProcessedEvent-10QL7FW4JP5V	Topic owner 577417548456								
Type <b>Standard</b>									
<a href="#">Subscriptions</a>   <a href="#">Access policy</a>   <a href="#">Delivery retry policy (HTTP/S)</a>   <a href="#">Delivery status logging</a>   <a href="#">Encryption</a>   <a href="#">Tags</a>									
<h3>Subscriptions (1)</h3> <table border="1"> <thead> <tr> <th>ID</th> <th>Endpoint</th> <th>Status</th> <th>Protocol</th> </tr> </thead> <tbody> <tr> <td>7d96445b-a4cb-4261-8021-f3c8b01dc27d</td> <td>sruathi.mallarapu@sjsu.edu</td> <td>Confirmed</td> <td>EMAIL</td> </tr> </tbody> </table>		ID	Endpoint	Status	Protocol	7d96445b-a4cb-4261-8021-f3c8b01dc27d	sruathi.mallarapu@sjsu.edu	Confirmed	EMAIL
ID	Endpoint	Status	Protocol						
7d96445b-a4cb-4261-8021-f3c8b01dc27d	sruathi.mallarapu@sjsu.edu	Confirmed	EMAIL						

## Successful Email Upon Completion of the Job

**Inbox** 393

Starred  
 Snoozed  
 Sent  
 Drafts 4  
 Chat +

No conversations Start a chat

AWS Notification Message [External](#) [Inbox](#)

AWS Notifications <no-reply@sns.amazonaws.com> Tue, May 18, 8:21 PM (5 days ago)

to me

{"version": "0.2", "id": "e2d31900-e638-556d-7b54-d83f035eab40", "detail-type": "Glue Job State Change", "source": "aws glue", "account": "577417548456", "time": "2021-05-19T03:21:35Z", "region": "us-east-1", "resources": ["arn:aws:glue:us-east-1:789:job/be223e19919379678c"], "detail": "Job run finished", "jobRunId": "d031ad46802a9e74dd09ea1f77a1c378b1d07789cab7be223e19919379678c", "message": "Job run succeeded!"}

--

If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:  
<https://sns.us-east-1.amazonaws.com/unsubscribeArn/arn:aws:sns:us-east-1:577417548456:demostack-SNSProcessedEvent-10C17FV4-JPSV-7d96445b-acb-4261-9027-0dbbd1c27d8&Endpoint=https://sns.mallorcaweb.s3.eu-central-1.amazonaws.com>

Please do not reply directly to this email. If you have any questions or comments regarding this email, please contact us at <https://aws.amazon.com/support>

## Redshift create table

The screenshot shows the Amazon Redshift Query Editor interface. At the top, there's a navigation bar with 'Amazon Redshift' and 'Query editor'. Below it is a toolbar with 'Editor', 'Query history', 'Saved queries', and 'Scheduled queries'. The main area has tabs for 'Query 2', 'Query 3', 'Query 4', 'Query 1', 'Query 5', and '+'. Each tab shows a list of SQL commands. The 'Query 1' tab is active and displays the following code:

```
1 CREATE TABLE public.dim_customer();
2 go
3 drop table dim_customer;
4 create table dim_customer(
5     customer_id int primary key,
6     first_name varchar(50),
7     last_name varchar(50),
8     email varchar(100),
9     phone_number varchar(20),
10    address_id int,
11    zip_code int,
12    city_id int,
13    country_id int,
14    gender char(1),
15    date_joined date,
16    last_update timestamp
17);
```

Below the code, there are buttons for 'Run', 'Save', 'Schedule', and 'Clear'. On the right side, there are status indicators for 'Status' (green), 'Connected' (green), 'database' (blue), 'dev' (blue), 'user' (blue), and 'invoker' (blue). A 'Change connection' button is also present. The bottom of the screen includes a feedback link and standard footer links for 'Feedback', 'English (US)', 'Privacy Policy', 'Terms of Use', and 'Cookie preferences'.

Amazon Redshift > Query editor

Editor    Query history    Saved queries    Scheduled queries

Query 1 +

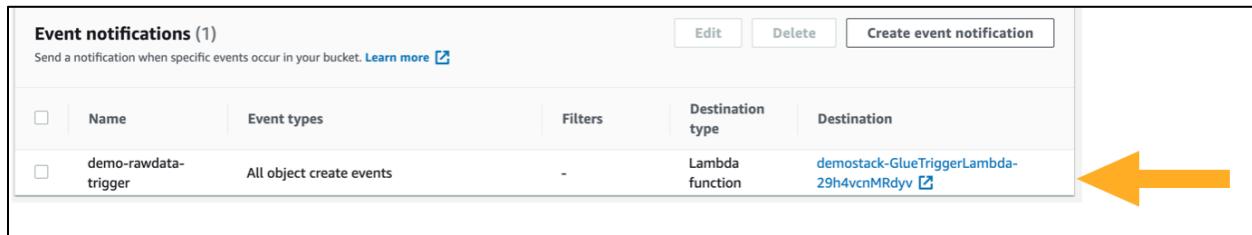
```
1 create table stations_1996_97(
2     stationid integer,
3     startdate varchar,
4     enddate varchar,
5     start_stations_id integer,
6     end_stations_id integer,
7     end_stations_name varchar,
8     end_stations_lat decimal,
9     end_stations_lng decimal,
10    average_rainfall,
11    total_rainfall,
12    percent_holiday,
13    percent_cold_start,
14    end_stations_latitude float,
15    end_stations_longitude float,
16    end_stations_imagine_float);
```

## Demo data flow

Created two S3 buckets namely “demo-rawdata” and “demo-processed-data”.  
“demo-rawdata” stores the raw data and “demo-processed-data” bucket stores the processed data.

<input type="radio"/>	<a href="#">demo-processed-data</a>	US East (N. Virginia) us-east-1	Objects can be public	May 17, 2021, 00:13:12 (UTC-07:00)
<input type="radio"/>	<a href="#">demo-rawdata</a>	US East (N. Virginia) us-east-1	Objects can be public	May 17, 2021, 00:12:47 (UTC-07:00)

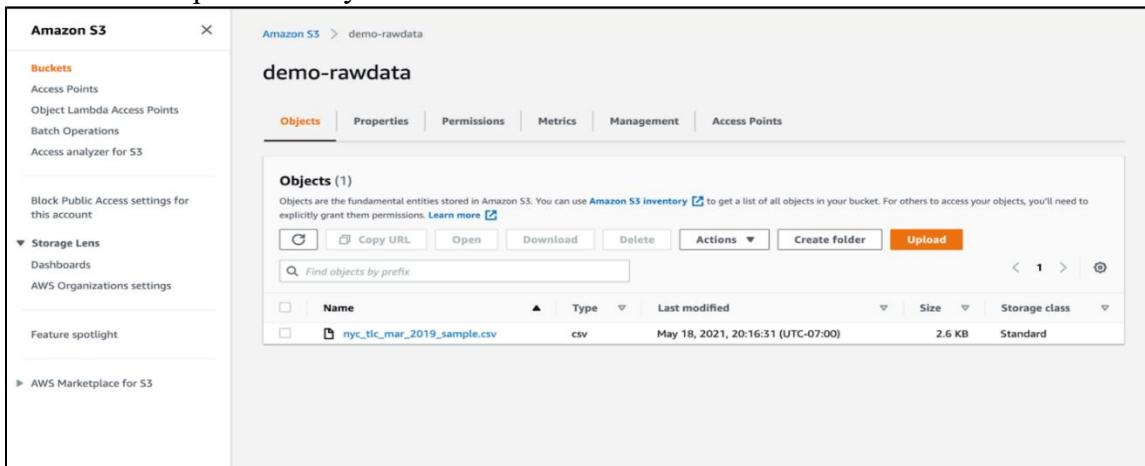
Created the S3 trigger name **demo-rawdata-trigger** in the”**demo-rawdata**” bucket. This trigger invokes the lambda function whenever new file is uploaded in the bucket.



**Event notifications (1)**  
Send a notification when specific events occur in your bucket. [Learn more](#)

Name	Event types	Filters	Destination type	Destination
demo-rawdata-trigger	All object create events	-	Lambda function	<a href="#">demostack-GlueTriggerLambda-29h4vcnMRdyv</a>

Uploaded the Sample Data of yellow taxi in the S3 bucket named “demo-rawdata”.



**Amazon S3** X

[Amazon S3](#) > demo-rawdata

### demo-rawdata

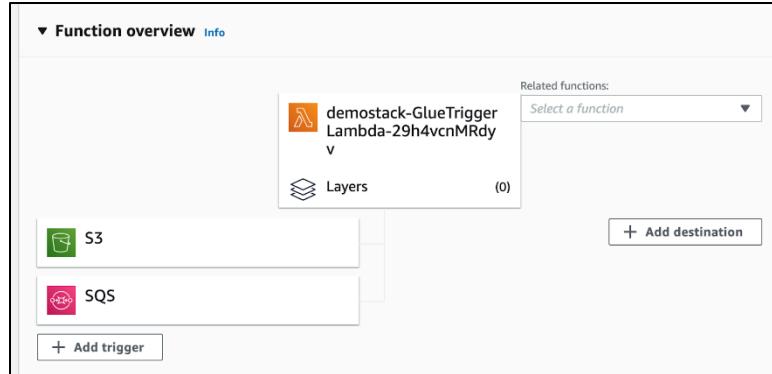
[Objects](#) [Properties](#) [Permissions](#) [Metrics](#) [Management](#) [Access Points](#)

**Objects (1)**

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](#)

Name	Type	Last modified	Size	Storage class
<a href="#">nyc_tlc_mar_2019_sample.csv</a>	csv	May 18, 2021, 20:16:31 (UTC-07:00)	2.6 KB	Standard

The Lambda function invoked above is used to create the crawler named “demo-crawler” and then creates the data catalog. SQS trigger is implemented for the retry logic.



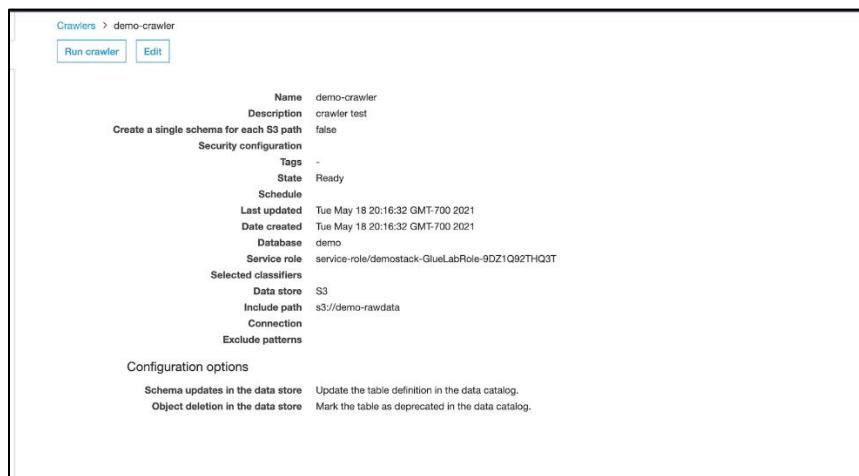
Code Snapshot for the above lambda function.

```

    }
} else{
    var dbName = 'demo';
    var params = {
        DatabaseInput: {
            Name: dbName,
            Description: 'Blog Post database',
        }
    };
    glue.createDatabase(params, function(err, data) {
        var params1 = {
            DatabaseName: dbName,
            Name: 'demo-crawler',
            Role: 'service-role/dl-data228-gluerole',
            Targets: {
                S3Targets: [{ Path: 's3://demo-rawdata' }]
            },
            Description: 'crawler test'
        };
        glue.createCrawler(params1, function(err1, data1) {
            startCrawler('demo-crawler', function(err2,data2){
                if(err2) callback(err2)
                else callback(null,data2)
            })
        });
    });
}
}

```

demo-crawler is created by the invocation of the above lambda function.



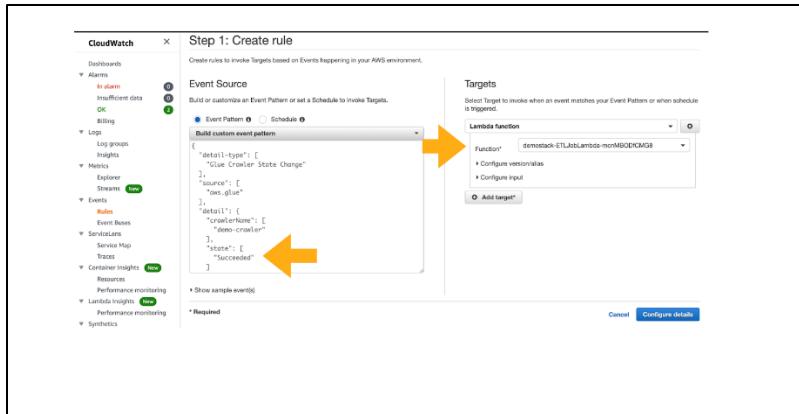
After the demo crawler crawls the data "demo-rawdata" catalog table is created.

Tables > demo_rawdata		Last updated 18 May 2021 08:17 PM	Table	Version (Current version)	▼
		<a href="#">View properties</a> <a href="#">Compare versions</a> <a href="#">Edit schema</a>			
<b>Name</b> demo_rawdata					
<b>Description</b>	demo				
<b>Database</b>	demo				
<b>Classification</b>	csv				
<b>Location</b>	s3://demo-rawdata/				
<b>Connection</b>					
<b>Deprecated</b>	No				
<b>Last updated</b>	Tue May 18 20:17:41 GMT-700 2021				
<b>Input format</b>	org.apache.hadoop.mapred.TextInputFormat				
<b>Output format</b>	org.apache.hadoop.hive.io.HiveIgnoreKeyTextOutputFormat				
<b>Serde serialization lib</b>	org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe				
<b>Serde parameters</b>	field.delim : , skip.header.line.count : 1 sizeKey : 2682 objectCount : 1 UPDATED_BY_CRAWLER : demo-crawler				
<b>Table properties</b>					
CrawlerSchemaVersion	1.0	recordCount	20	averageRecordSize	131
CrawlerSchemaDeserializerVersion	1.0	compressionType	none	columnsOrdered	true
areColumnsQuoted	false	delimiter	,	typeOfData	file

## Schema of the raw data

Column name	Data type	Partition key
1 vendorid	bigint	
2 tpep_pickup_datetime	string	
3 tpep_dropoff_datetime	string	
4 passenger_count	bigint	
5 trip_distance	double	
6 ratecodeid	bigint	
7 store_and_fwd_flag	string	
8 pulocationid	bigint	
9 dolocationid	bigint	
10 payment_type	bigint	
11 fare_amount	double	
12 extra	double	
13 mta_tax	double	
14 tip_amount	double	
15 tolls_amount	double	
16 improvement_surcharge	double	
17 total_amount	double	
18 congestion_surcharge	double	

An event rule is created in Cloud Watch that invokes another lambda function upon successful crawling of the data.



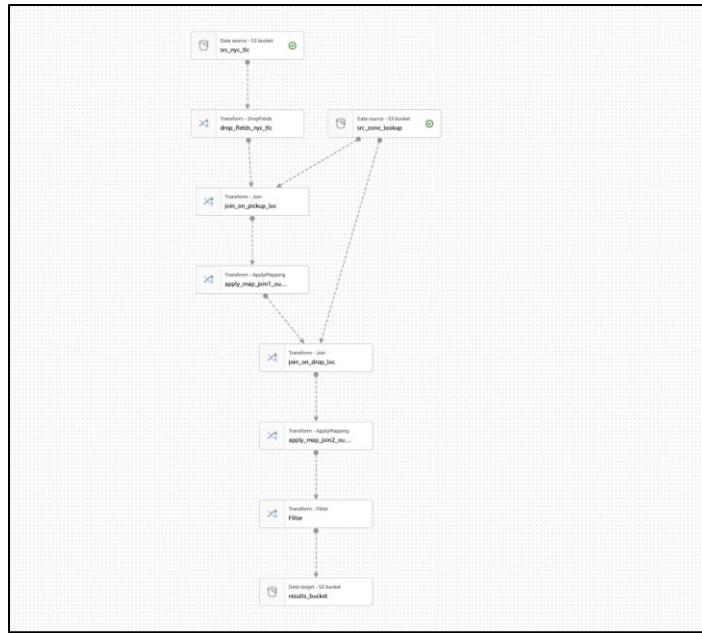
Code snapshot for the second lambda function. It invoke the ETL job named 'demo-job'

```

1 var AWS = require('aws-sdk');
2 var sns = new AWS.SNS({ region: "us-east-1" });
3 var S3 = new AWS.S3();
4 var glue = new AWS.Glue({apiVersion: '2017-03-31'});
5 exports.handler = function(event, context, callback) {
6     var params = {
7         JobName: 'demo-job',
8         Timeout: 20,
9     };
10    glue.startJobRun(params, function(err1, data1) {
11        if (err1) {
12            console.log(err1, err1.stack);
13        } else {
14            console.log(data1);
15        }
16    });
17    console.log(JSON.stringify(event, null, 3));
18};

```

The ETL job is created using AWS Glue studio that cleans the data and transfers it to processed S3 bucket. Cleaning process involved dropping unnecessary columns, changing data types, filtering the data, and performing the joins with mapping files.



Processed Data is stored in S3 processed buckets.

Node properties	Data target properties - S3
Output schema	<input type="text"/>
Format	CSV
Compression Type	None
S3 Target Location	Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/). <input type="text" value="s3://demo-processed-data/"/> <input type="button" value="View Details"/>
	<input type="button" value="Browse S3"/>
<input type="checkbox"/> Data Catalog update options <a href="#">Info</a> Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3 backed table. <input type="radio"/> Do not update the Data Catalog <input checked="" type="radio"/> Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions <input type="radio"/> Create a table in the Data Catalog and on subsequent runs, existing schema and add new partitions	

The data is stored in processed S3 bucket.

Created the S3 trigger name **trigger-gluejob** in the "**demo-processed-data**" bucket. This trigger invokes the lambda function as soon as the processed data comes.

This Lambda function invokes "**demo-redshift-job**"

```

1 var AWS = require('aws-sdk');
2 var sns = new AWS.SNS({ region: "us-east-1" });
3 var s3 = new AWS.S3();
4 var glue = new AWS.Glue({apiVersion: '2017-03-31'});
5 exports.handler = function(event, context, callback) {
6   var params = {
7     JobName: 'demo-redshift-job',
8     Timeout: 20,
9   };
10  glue.startJobRun(params, function(err1, data1) {
11    if (err1) {
12      console.log(err1, err1.stack);
13    } else {
14      console.log(data1);
15    }
16    console.log(JSON.stringify(event, null, 3));
17  });
18}

```

Configured the redshift Connection that is required by the **demo-redshift-job** in order to transfer the data from processed S3 bucket to Redshift.

Successful execution of **demo-redshift-job**. It means data has transferred to Redshift tables

The screenshot shows the AWS Glue Job History page for the job 'demo-redshift-job'. The job is listed as 'Spark python' with an ARN of 's3://aws-glue...'. It was run on 17 May 2021 at 12:23 PM and is currently 'Disable'. The 'Metrics' tab is selected. A table below shows the run details:

Run ID	Retry attempt	Run status	Error	Output Logs	Error logs	Glue version	Maximum capacity	Triggered by	Start time	End time	Start-up time	Execution time	Timeout	Delay	Job run input
jr_dd031ad4...	-	Succeeded		Logs	Error logs	2.0	10		18 ...	18 ...	18 secs	1 min	20 mins		s3://aws-glu...

Another event rule is created in Cloud Watch that invokes SNS services. This service notifies through an email upon successful completion of the **demo-redshift-job**.

The screenshot shows the AWS CloudWatch Events Rule configuration for the rule 'demostack-OpsEventRule-1B9RLHBNUEMRM'. The 'Event pattern' section contains the following JSON code, with an orange arrow pointing to the 'state': 'SUCCEEDED' part:

```
{
  "detail-type": "Glue Job State Change",
  "source": [
    "aws.glue"
  ],
  "detail": {
    "jobName": [
      "demo-redshift-job"
    ],
    "state": [
      "SUCCEEDED"
    ]
  }
}
```

The 'Targets' section shows a single target:

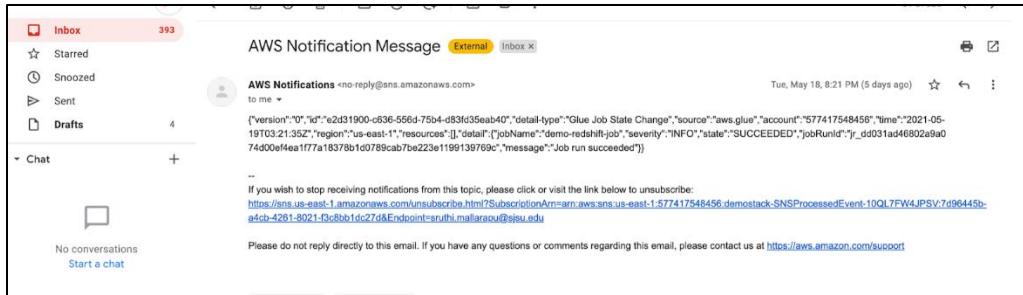
Type	Name	Input	Role	Additional parameters
SNS topic	demostack-SNSProcessedEvent-10QL7FW4JPSV	Matched event		

Configured Email in AWS SNS.

The screenshot shows the AWS SNS Topic Details page for the topic 'demostack-SNSProcessedEvent-10QL7FW4JPSV'. The 'Subscriptions' tab is selected, showing one subscription:

ID	Endpoint	Status	Protocol
7d96445b-a4cb-4261-8021-f3c8bb1dc27d	sruthi.mallarapu@sjtu.edu	Confirmed	EMAIL

Successful Email Upon Completion of the Job



## Data analysis in Redshift

We stored our processed data as tables in Redshift and executed SQL queries for analysis.

As a first step, a cluster must be created to host the database instance. Subsequently, a database is created with name and password and becomes functional.

### Redshift Cluster

General information			
Cluster identifier redshift-cluster-1	Status <span style="color: green;">Available</span>	Node type dc2.large	Endpoint <a href="#">redshift-cluster-1.cqk8mnrhrted.us-east-1.amazonaws.com</a>
Cluster namespace 22f8a8f9-7a01-473e-b73d-0c42e6d6626d	Date created May 15, 2021, 03:54(UTC-07:00)	Number of nodes 1	JDBC URL <a href="#">jdbc:redshift://redshift-cluster-1.cqk8...</a>
	Storage used 3.16% (5.06 of 160 GB used)	AQUA Not available	ODBC URL <a href="#">Driver={Amazon Redshift (x64)}; Serve...</a>

**Properties**

Database configurations			
Database name dev	Parameter group Defines database parameter and query queues for all the databases. <a href="#">default.redshift-1.0</a>	Encryption Disabled	Audit logging Disabled
Port 5439	SSH ingestion setting (cluster public key) <a href="#">ssh-rsa AAAAB3NzaC1yc2EAAAQ...</a>	AWS KMS key ID -	
Admin user name awsuser			

Figure 4: Cluster created along with the database.

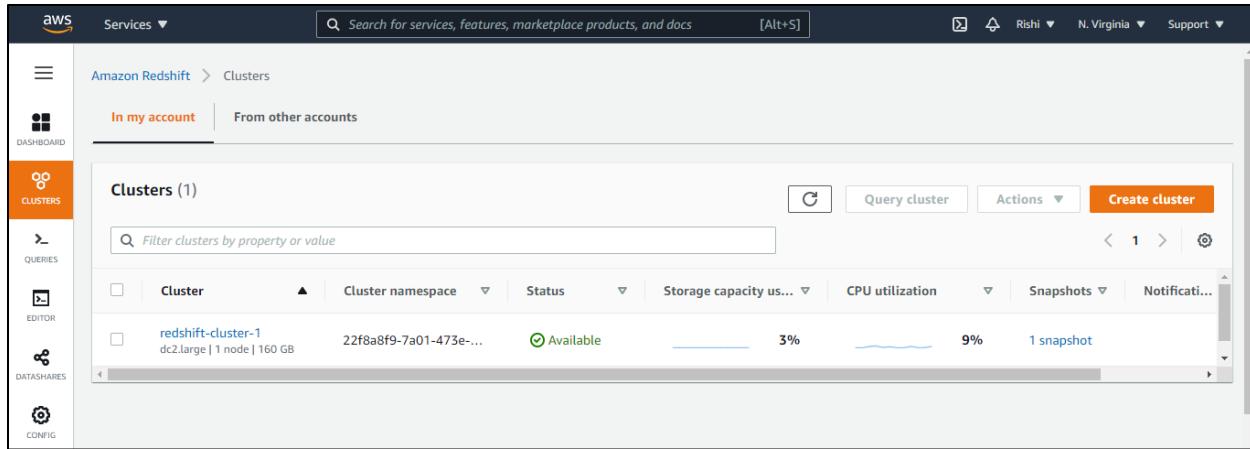


Figure 5: Redshift cluster is in status Available, ready to query data for analysis.

## Our Taxi and Citi bike data tables in Redshift

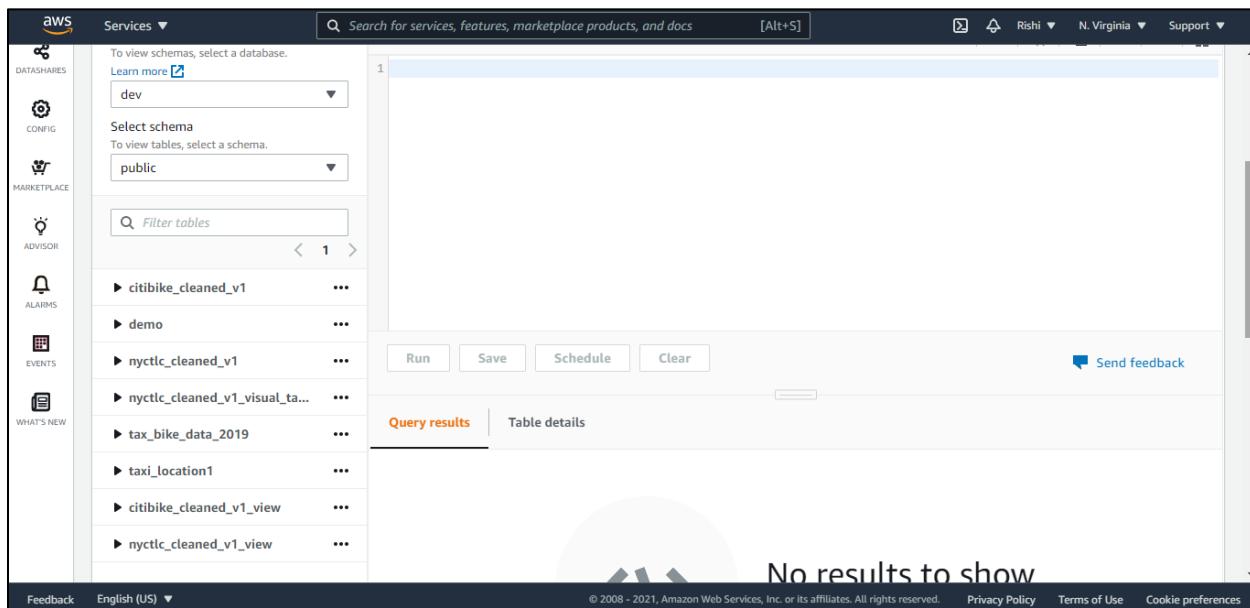


Figure 6: Redshift data tables

## Queries, loads and system performance parameters are tracked automatically by AWS

The screenshot shows the AWS CloudWatch Metrics Queries and loads interface. At the top, there's a search bar and navigation links for Rishi, N. Virginia, and Support. Below the header, a table lists 68 queries. The columns include Start time, Query ID, Status, Duration, and SQL. Most queries are completed and took between 1 min and 24 sec. The SQL code for each query is partially visible, showing various QuickSight-related operations.

Figure 7: Infrastructure managed by the cloud vendor

## Some analysis using SQL:

The screenshot shows the Amazon QuickSight SQL editor. The query is as follows:

```

5 -- Get the average distance and average cost per mile by do_zone
6 SELECT do_zone,
7       avg(trip_distance) AS avgDist,
8       avg(total_amount/trip_distance) AS avgCostPerMile
9 FROM nyccta_cleansed_v1
10 WHERE trip_distance > 0
11 AND total_amount > 0
12 GROUP BY do_zone;
13
  
```

Below the query, there are buttons for Run, Save, Schedule, Clear, and Send feedback. The results section shows the query was completed on May 18, 2021, at 14:41:32, with an elapsed time of 00 m 02 s. The results table has columns do\_zone, avgdist, and avgcostpermile, containing data for various New York City zones.

Figure 8: SQL Analysis - Taxi data: Average distance, average cost per mile

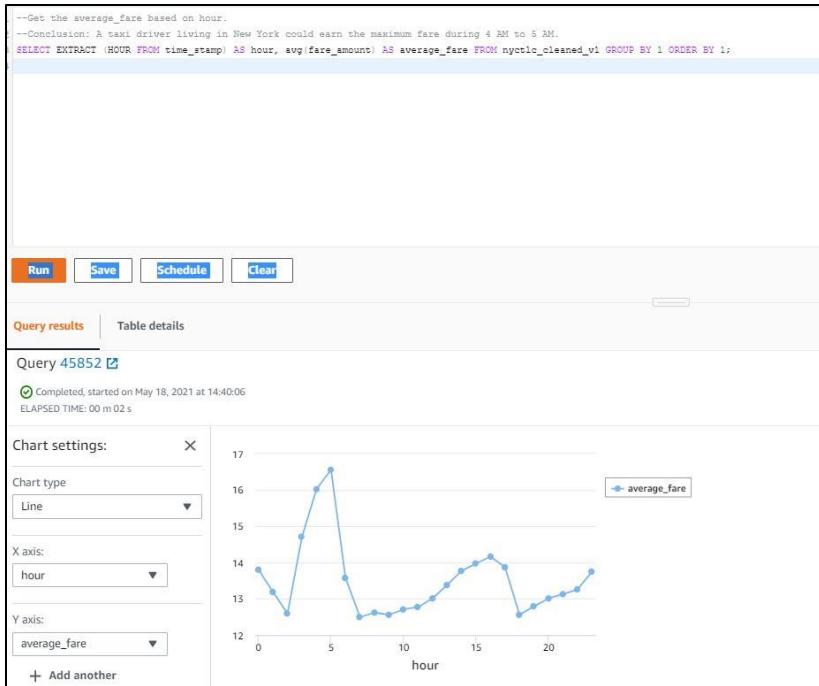


Figure 9: SQL Analysis - Taxi data: Average fare based on hour.

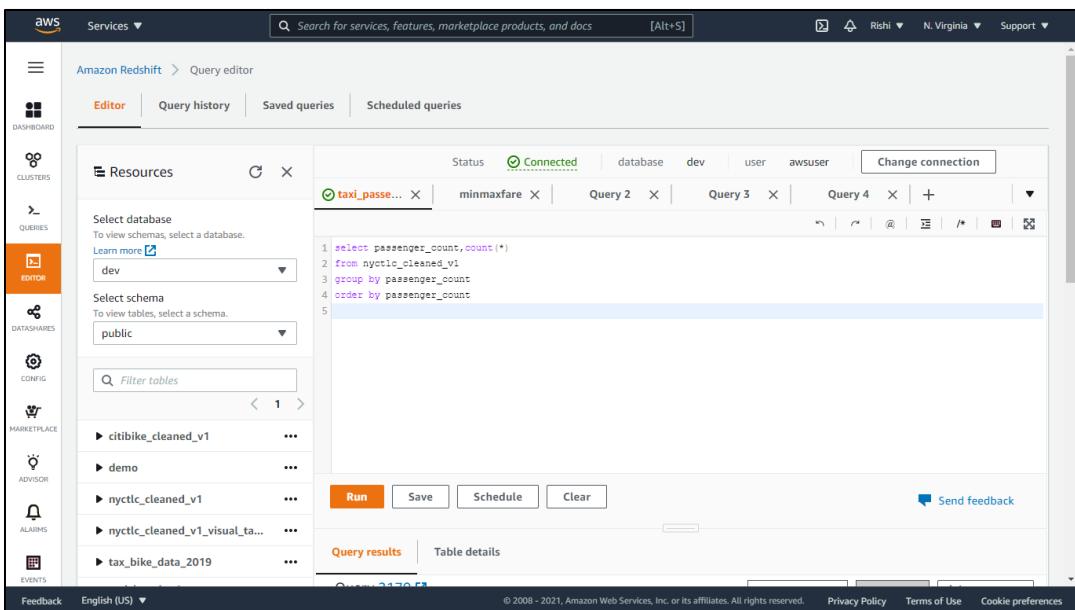


Figure 10: SQL Query: Taxi data - Passenger count

The screenshot shows the AWS Redshift Query results interface. On the left, there's a sidebar with 'EVENTS' and 'WHAT'S NEW'. The main area displays a table titled 'passenger\_count' with the following data:

passenger_count	count
0	499949
1	18546531
2	3895347
3	1084913
4	502578
5	1074978
6	642639
7	110
8	79
9	70

At the bottom, there are links for 'Feedback', 'English (US)', and 'Cookie preferences'.

Figure 11: SQL Query result - Passenger count

## Data Visualizations in QuickSight

Data has been imported from Redshift to QuickSight.

### Connection process from Redshift to QuickSight:

Create an IAM role for Redshift to QuickSight and add the same role in VPC Security Group at Redshift.

### Datasets

The below datasets have been imported from AWS Redshift and will be used to generate our graphs and charts.

The screenshot shows the AWS QuickSight Datasets page. On the left, there's a sidebar with 'Favorites', 'Recent', 'Dashboards', 'Analyses', and 'Datasets'. The main area displays a table of datasets:

Name	Owner	Last Modified
tax_bike_data_2019	Me	5 days ago
nyctlc_cleaned_v1	Me	7 days ago

At the top right, there's a 'New dataset' button.

Figure 12: Datasets with NYC yellow Taxi and Citi Bike data for visualizations

### Analyses and Dashboard

From the above dataset, analysis was performed with the required visual type. We have options for filtering and performing any required action with the data.

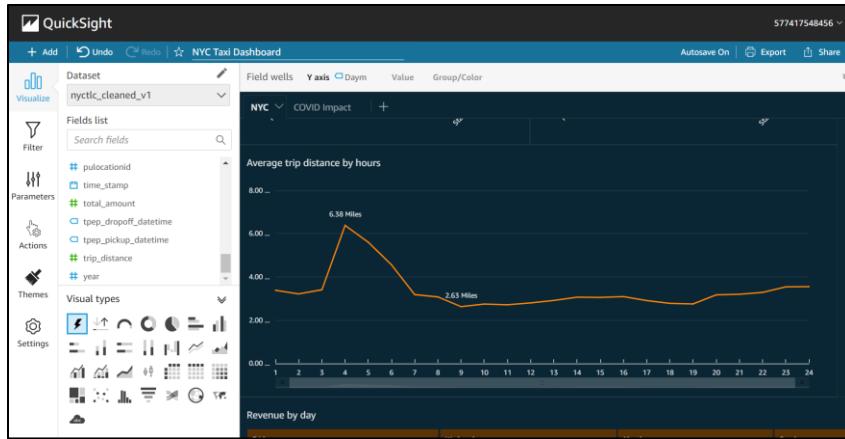


Figure 13: Analysis window with menu options

Once the graph and visualization are finalized, the same can be published as Dashboards using the ‘share’ option. The published dashboards can be viewed under Dashboards menu in QuickSight.

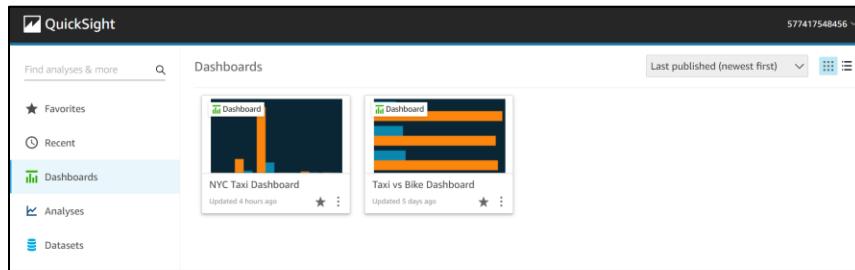
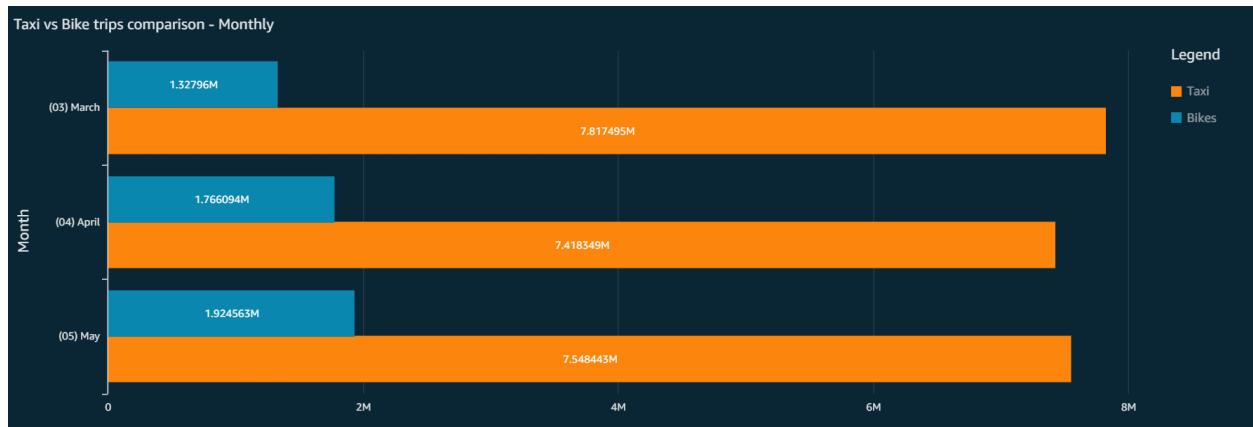


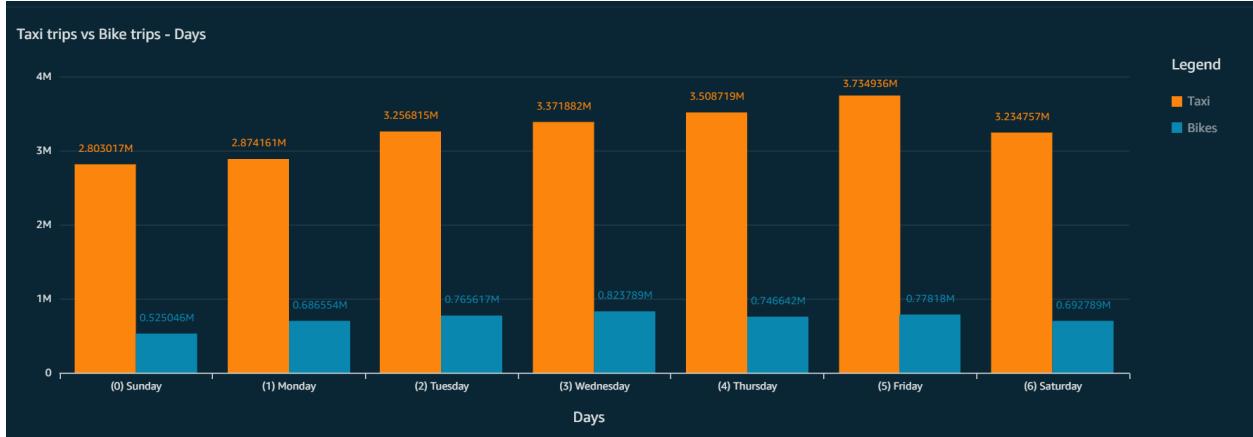
Figure 14: Dashboards menu

### Visualization for NYC – Yellow Taxi and Citi Bike:

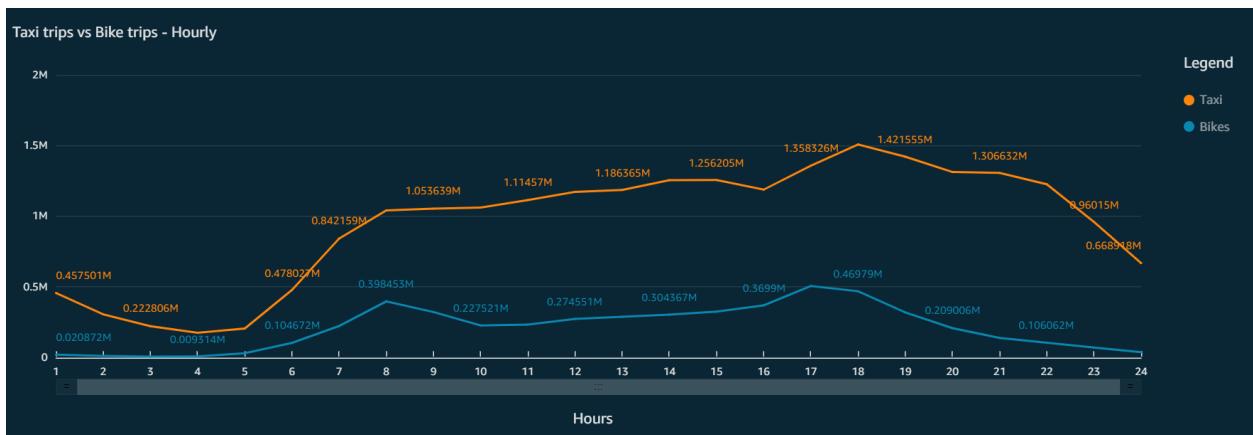
The below bar graph represents the count of trips comparison between the Taxi and Citi Bike. Here, we could see the increase in bike trips as month goes on.



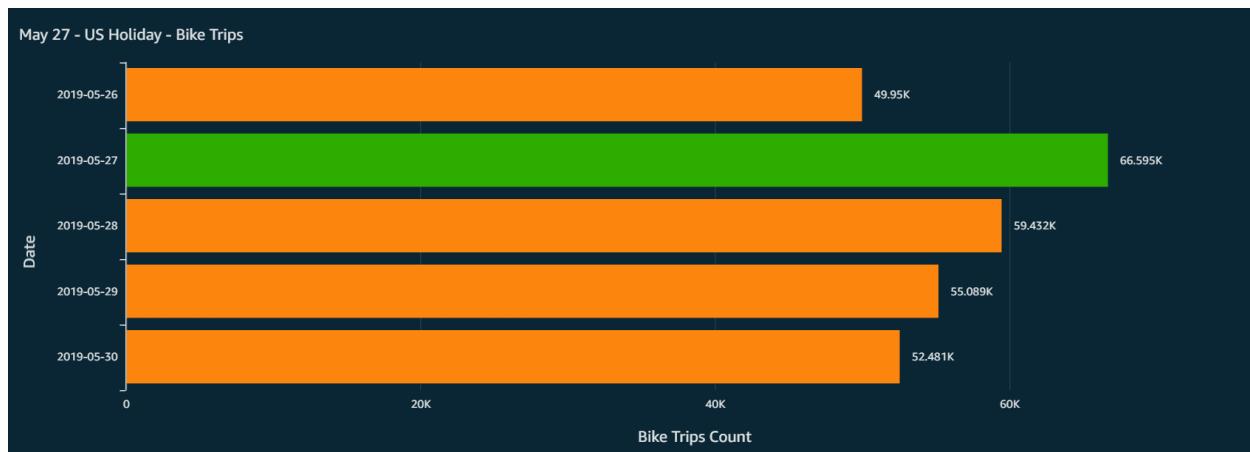
In the below bar graph, we have segregated the trips of both Taxi and Bike based on the day. We have analyzed that, peak day for taxi is Friday and bike is Wednesday. The lowest trips for Taxi and Bike are on Sunday.



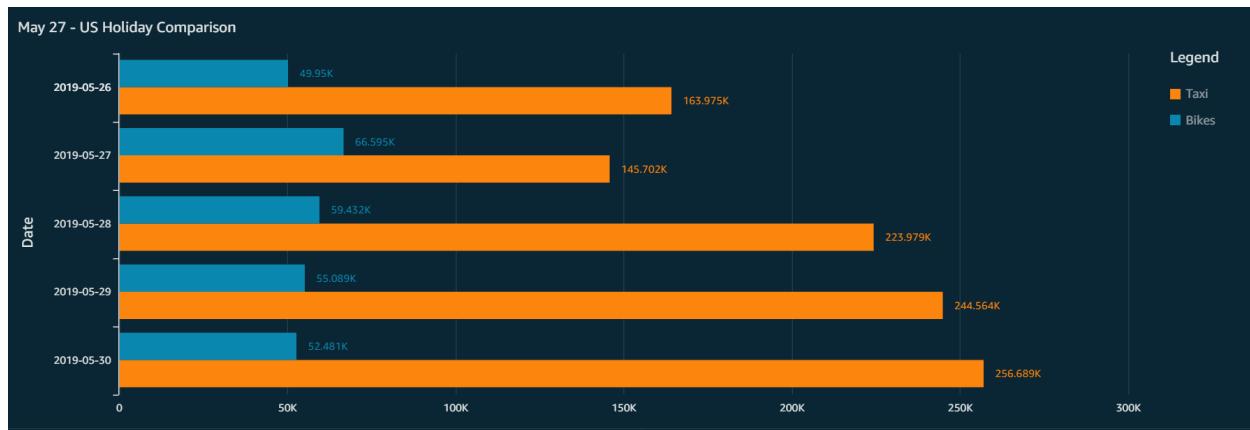
The below line graph implies the peak booking count based on hours and we analyzed that 8 am is the highest on morning and 5 pm for Bike and 6 pm for taxi is the highest trips witnessed hours of the day. This time is claimed to be the office and school's closing time.



We made analysis for May 27, 2019 from our dataset, which was a public holiday, and we could see, on that particular day, the Bike trips is comparatively high on the week.



The same public holiday trips count has been analyzed for both Taxi and Bike in a single graph and the outcome was negative to Taxi.



From our dataset, we have found that the average minutes for both bike and taxi to increase the availability.



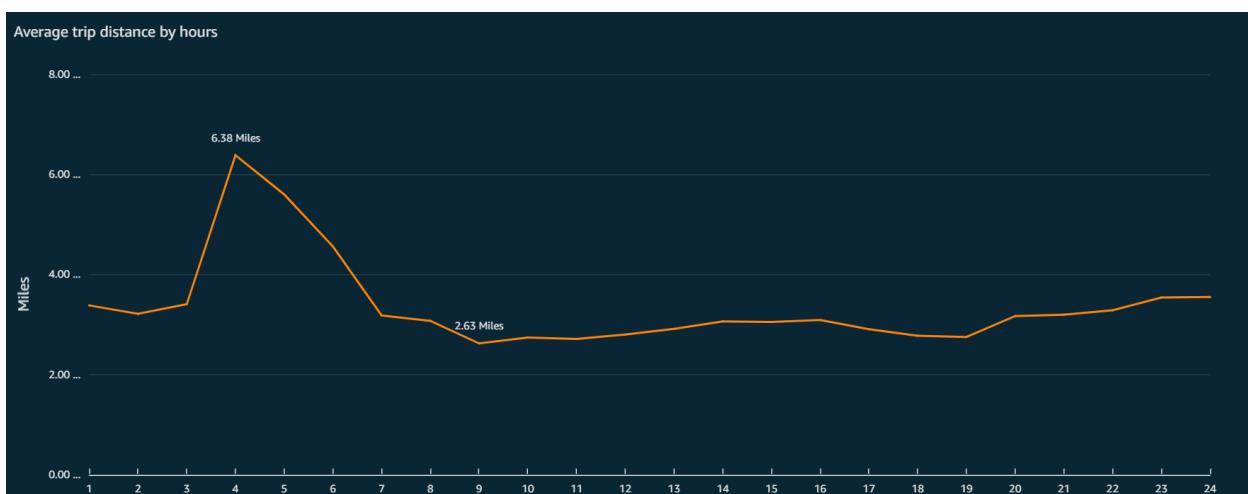
### Visualization for NYC yellow Taxi and COVID-19 impact:

As yellow taxi generates significant revenue for NYC, we were additionally inspired to perform some analysis on the same and the effects of COVID on the business.

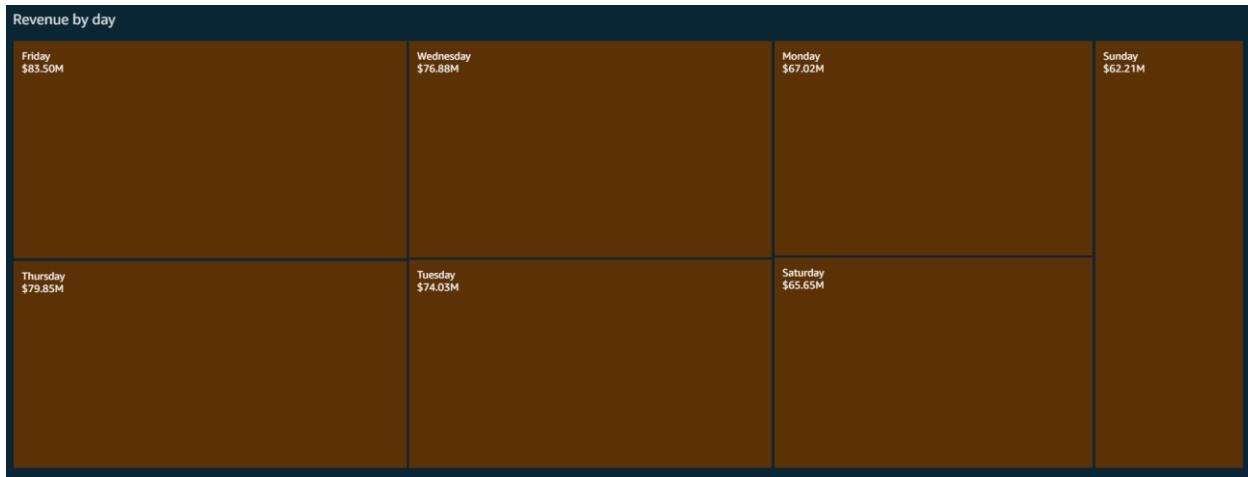
The below graph represents the trip count and revenue of top 6 boroughs. Manhattan City holds the highest trip count and the highest revenue.



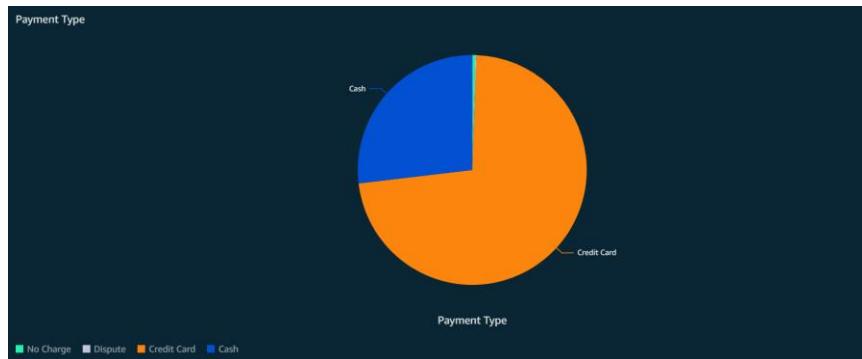
This graph provides information about the average distance (Miles) travelled based on hours and the result is that 4 am has the highest with average of 6.38 Miles and the lowest occurs at 9 am with 2.63 Miles.



The tree map below shows the revenue generated by NYC on the each day of the week. Friday, being the start of the weekend holds the highest revenue compared to other days.

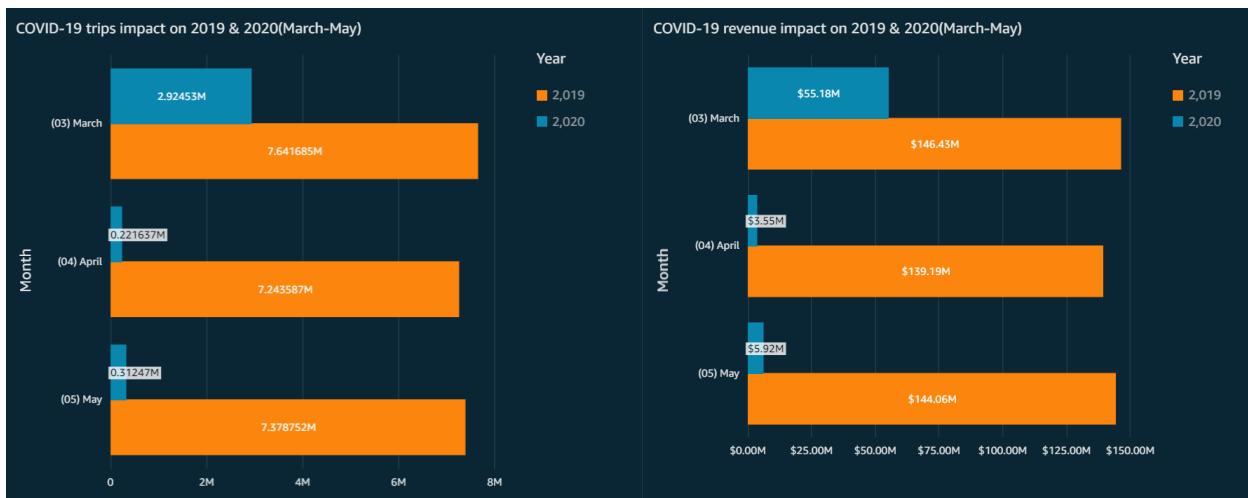


The pie chart below represents the payment types used by the customers.

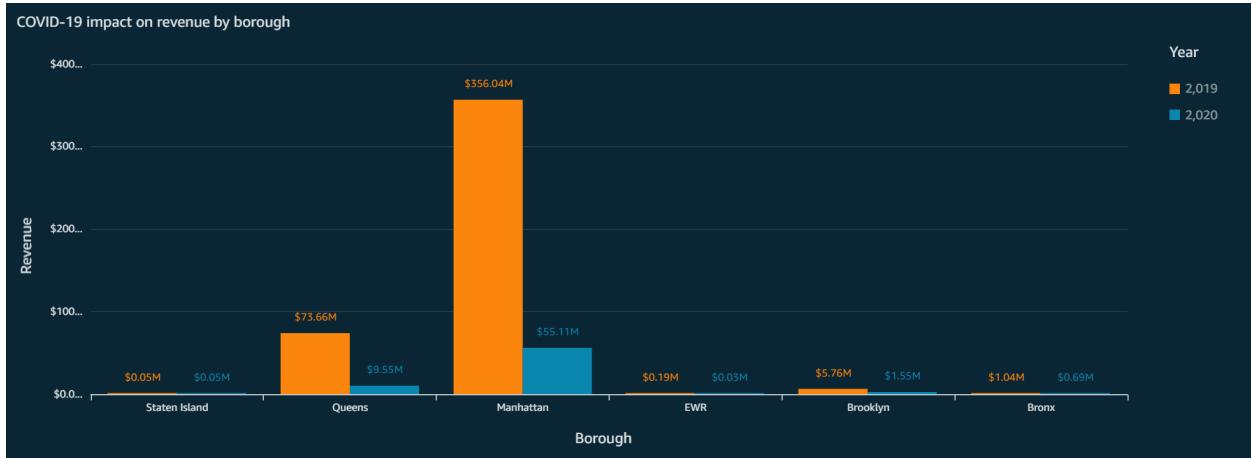


### COVID-19 Impact:

The bar graph below clearly shows the impact of the global pandemic and shows the trips and revenue impact of the month of March – May in 2019 and 2020.



The revenue impact on the business due to COVID-19 is analyzed for the top 6 boroughs on the below bar graph.



## **VI. Team Members and Roles**

As a group, we invested some time in fixing on the goals of this project and brainstorming ideas. It involved discussions from choosing the dataset, coming up with relevant use-cases and to finalizing the tools we will be using to implement the project.

Once fixed on our datasets and the tools, we tried to break up our project tasks into parts and assigned ourselves a significant portion of the project. For example, one of us proceeded to process the raw data, while the other tried to learn about visualization. In this way, many different components of our project were simultaneously developed.

From time to time, we worked together by means of recurring meetings and weekly sprints. During the meetings, everyone shared their updates along with challenges faced and the entire team would brainstorm to solve those challenges. As a result, the entire team was always updated and involved in all parts of the project.

## **VII. Future Scope**

We believe this project can be improved further by adding even more data or processes. One such endeavor we took up was to combine real-time data with the data we already had.

We wanted to know how useful the Twitter data of NYC TLC is to the public and the business. So we pulled real-time Twitter data using node.js scripts and loaded our S3 buckets with the tweets using Amazon Kinesis Data Firehose delivery stream.

The below screenshot shows the Kinesis data stream 'Twitterstream228', created to push real-time tweets to S3 bucket named 'fhbucket228'.

The screenshot shows the AWS Kinesis Data Firehose delivery streams interface. The left sidebar has a 'Data Firehose' section selected. The main area displays a table with one row, showing a delivery stream named 'Twitterstream228'.

Name	Status	Creation time	Source	Data transformation	Destination
Twitterstream228	Active	2021-05-17T14:18-0700	Direct PUT and other sources	Disabled	Amazon S3 fbucket228

The below screenshot shows the twitter data being pulled using Twitter API by running node.js scripts.

The below screenshot shows the twitter data getting created in S3 in real-time.

Objects (5)									
Objects are the fundamental entities stored in Amazon S3. You can use <a href="#">Amazon S3 inventory</a> to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. <a href="#">Learn more</a>									
<input type="checkbox"/>	<a href="#">Copy URL</a> <a href="#">Open</a> <a href="#">Download</a> <a href="#">Delete</a> <a href="#">Actions ▾</a> <a href="#">Create folder</a> <a href="#" style="background-color: orange; color: white; padding: 2px 10px;">Upload</a>								
<input type="text"/> <a href="#">Find objects by prefix</a>									
Name	▲	Type	▼	Last modified	▼	Size	▼	Storage class	▼
<input type="checkbox"/> <a href="#">Twitterstream228-1-2021-05-18-04-30-12-1def2318-78a2-47d6-8d5-90f780ed226b</a>	-			May 17, 2021, 21:35:14 (UTC-07:00)		29.7 KB		Standard	
<input type="checkbox"/> <a href="#">Twitterstream228-1-2021-05-18-04-36-02-74dd0169e-bdfa-42a1-9179-2567e3e98f6a</a>	-			May 17, 2021, 21:41:04 (UTC-07:00)		9.4 KB		Standard	
<input type="checkbox"/> <a href="#">Twitterstream228-1-2021-05-18-04-42-09-2039c2c8-6f16-40fd-a01f-f955e051ab01</a>	-			May 17, 2021, 21:47:11 (UTC-07:00)		23.6 KB		Standard	
<input type="checkbox"/> <a href="#">Twitterstream228-1-2021-05-18-04-48-24-b679c407-4f33-4fd3-ab8d-b9347a3dddbe</a>	-			May 17, 2021, 21:53:26 (UTC-07:00)		22.7 KB		Standard	
<input type="checkbox"/> <a href="#">Twitterstream228-1-2021-05-18-04-54-12-4135ec63-efeb-4afc-e4f926f4c80</a>	-			May 17, 2021, 21:59:14 (UTC-07:00)		3.8 KB		Standard	

This twitter data can be loaded into Redshift or ‘Elasticsearch’ and Kibana for data analysis. This data analysis can be also integrated with our data for comparison and further analysis.

Additionally, we also converted shapefile available in the taxi data repository to geojson using online software convertors and loaded this data into DynamoDB for visualization. Unfortunately, the conversion was flawed and the latitude, longitude values were not helpful. We are planning to debug this more and look for alternative methods to extract location field values from shapefiles.

```
geojson_name=r'C:\Users\vijay\Documents\A\SJSU\01.BigData\dynamodb_table_name='taxi1')
print(f"{NUM_FEATURES_UPLOADED} features in data uploaded to dynamodb")

Connecting to Dynamodb...
Enter your aws_access_key:
AKIAYM4GESKUNZEWT5UL
Enter your aws_secret_access_key:
HnRDOZCM8z0EJv4+uEv4Nb1WxJA90+mxenbmzjW
Creating new table...
Uploading segments to table...
263 features in data uploaded to dynamodb
```

Figure 15: Python script uploaded data to DynamoDB

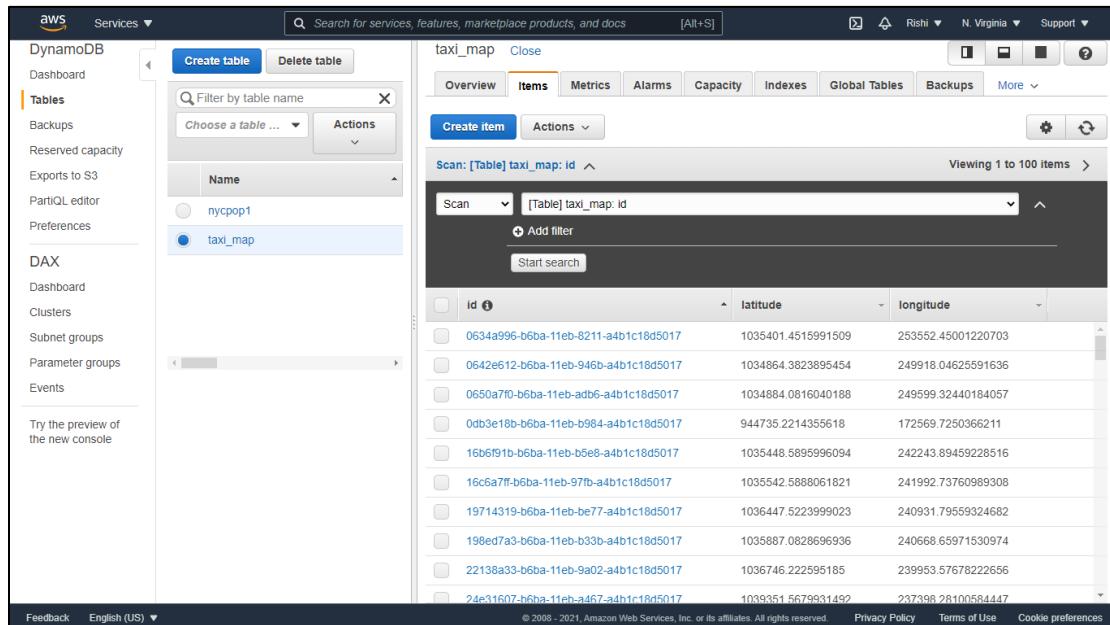


Figure 16: Latitude, longitude data is invalid.

## VIII. Learnings and Conclusion

We were able to achieve our goals by leveraging cloud services provided by AWS. We worked on big data and applied analysis on big data. As a result, this project helped us significantly in

- Exploring plethora of AWS services like automated ETL jobs that saves the time and manual effort and applies to huge datasets.
- Uncovering issues and debug errors related to cloud technology and software. IAM is one aspect that was totally new and was required to be properly configured for every step to work.

- Reading, understanding and interpreting logs - AWS creates logs for every action using its in-built service ‘CloudWatch’ which is useful in debugging errors with any of the services.
- Uncovering issues with big data and set those right. Loading, processing, debugging issues with big data took time and made us realize the significance of cloud infrastructures in today’s expansive data world.
- Exploring analytics by performing querying and join operations on tables with millions of rows made possible due to fast, fully managed, petabyte scale Redshift.
- Bringing about visualizations in the form of graphs and charts in AWS QuickSight.

By the end of the project, we were able to learn a lot about cloud technology and AWS and we attribute our successful completion to our professor and ISA. We sincerely thank the professor for this opportunity and his every lecture and guidance. We also thank our ISA for continued and detailed guidance during every stage of the project.

## **IX. References**

- <https://www1.nyc.gov/site/tlc/index.page>
- <https://zacks.one/aws-athena-lab/>
- <https://docs.aws.amazon.com/ses/latest/DeveloperGuide/event-publishing-redshift-firehose-stream.html>
- NYC TLC data source : [https://s3.amazonaws.com/nyc-tlc/trip+data/yellow\\_tripdata\\_2020-01.csv](https://s3.amazonaws.com/nyc-tlc/trip+data/yellow_tripdata_2020-01.csv)
- Citi bike data: <https://s3.amazonaws.com/tripdata/index.html>
- [https://docs.aws.amazon.com/code-samples/latest/catalog/code-catalog-python-example\\_code-kinesis.html](https://docs.aws.amazon.com/code-samples/latest/catalog/code-catalog-python-example_code-kinesis.html)
- <https://docs.aws.amazon.com/quicksight/latest/user/enabling-access-redshift.html>