**DAMG 7290 Data Warehousing & Bus Intel SEC 01**
**Spring 2022**

**Team-4**
Project Portfolio
**Aviation Data Warehousing**

**Members**
Siddhesh Aher
Bhavya Batra
Ronen Sengupta
Aishwarya Melige Suryanarayana

# Table of Contents

# 1.0 Introduction

Aviation data is available everywhere but tracking the airport and airline history physically is not easy. The objective of this project is to build a Data Warehouse using airports, airlines, routes, and passenger datasets taken from various sources. The Data transformations will be done in the staging area and visualizations will be created using Tableau/Power BI to get some more general insights.

The seven datasets considered for this project are taken from various sources. These datasets provide information about airports, airlines operated in different countries along with their respective IATA, ICAO codes, and routes between two airports. It also gives information on the number of passengers by different airlines.

The Aviation Data Warehouse will not only give historical insights but also provides an important means to effectively support airport decision-making. It will allow users to access critical data from several sources in a single place. It can be used to understand traffic patterns of airlines for different years, track aircraft history, observe airport congestion, etc.

## 1.0.1 Business Scenario

To gather appropriate data, we start the process by identifying the input sources. We consider the flight dataset obtained from Kaggle and Flight.org as the main source of input. We then extract the data present in the CSV format and store it in a staging area. This will help in resolving inaccuracies and errors in the data and verifying the input based on the required conditions.

To showcase the trends and extract significant information from the Data warehouse we outline a dimension model containing measures describing the Number of Passengers and Airline Routes through attributes like Passenger_Count, Source Airport, Destination Airport.

In order to analyze the changes in the dataset, we incorporate Slowly Changing Dimensions on a few attributes present in the dimension tables to keep track of the variation and history of the data. This will help us analyze the data more precisely.

Here are a few scenarios we considered to understand the slowly changing dimensions.

- Keeping track of the changes in Flight and the Number of Airlines will help us understand the availability of the Airline and flight modification.
- Keeping track of the Passengers around the airport around the year can help us with planning proper facilities for the passenger.
- Keeping a track of the number of Airports in the Country can help us understand the Passenger_country over a particular period.

Through the OLAP cube analyzing this data with historical values can help us provide valuable information in order to improve the performance of the Flights. To display the visualizations, we create dashboards containing all the crucial information via Tableau/Power BI.

## 2.0 Data Description

The datasets are considered from Flights.org and Kaggle. The project has seven datasets. The description of the datasets are as follows:

AIRPORT DATASET: This Dataset gives information about airports spanning the globe. This dataset has a total of 14 attributes and 7699 datapoints.

| Attributes | Data Types | Description |
|------------|------------|-------------|
| Airport_ID | Int | Unique identifier for the Airport. Primary Key(PK) |
| Name | Char | Name of airport |
| City | Char | Main city served by airport. May be spelled differently from Name |
| Country | Char | Country or territory where airport is located |
| IATA | Char | 3-letter IATA code. Null if not assigned/unknown |
| ICAO | Char | 4-letter ICAO code. Null if not assigned |
| Latitude | Float | Decimal degrees, usually to six significant digits. Negative is South, positive is North |
| Longitude | Float | Decimal degrees, usually to six significant digits. Negative is West, positive is East |
| Altitude | Int | In feet |
| Time zone | Int | Hours offset from UTC. Fractional hours are expressed as decimals, e.g., India is 5.5 |

| | | |
|---|---|---|
| DST | Char | Daylight savings time. One of E (Europe), A (US/Canada), S (South America), O (Australia), Z (New Zealand), N (None) or U (Unknown) |

AIRLINE DATASET: It gives information about the airlines operating around the globe along with their unique IATA/ICAO codes.

It has a total of 8 attributes and 6163 records.

| Attributes | Data Types | Description |
|---|---|---|
| Airline_ID | Int | Unique identifier for this airline. Primary Key(PK) |
| Name | Char | Name of the airline |
| IATA | Char | 2-letter IATA code, if available |
| ICAO | Char | 3-letter ICAO code, if available |
| Call sign | Char | Airline callsign |
| Country | Char | Country or territory where airport is located |
| Active | Boolean | "Y" if the airline is or has until recently been operational, "N" if it is defunct |

ROUTE DATASET: It gives the information on the routes between two airports and the number of stops for that route taken. It has 10 attributes and 67664 records.

| Attributes | Data Types | Description |
|---|---|---|
| Route_ID | Int | Primary Key(PK) |
| Airline | Char | 2-letter (IATA) or 3-letter (ICAO) code of the airline |
| Airline_ID | Int | Foreign Key (FK) |
| Source Airport | Char | 3-letter (IATA) or 4-letter (ICAO) code of the source airport |
| Source Airport_ID | Int | Unique identifier for source airport |

| Destination Airport | Char | 3-letter (IATA) or 4-letter (ICAO) code of the destination airport |
|---|---|---|
| Destination Airport_ID | Int | Unique identifier for destination airport |
| Codeshare | Boolean | "Y" if this flight is a codeshare (that is, not operated by Airline, but another carrier), empty otherwise |
| Stops | Int | Number of stops on this flight ("0" for direct) |
| Equipment | Char | 3-letter codes for plane type(s) generally used on this flight, separated by spaces |

COUNTRIES DATASET: It has all the countries listed along with their ISO codes. It has 3 attributes and 262 records.

| Attributes | Data Types | Description |
|---|---|---|
| Country_ID | Int | Primary Key (PK) |
| Name | Char | Full name of the country or territory |
| iso_code | Char | Unique two-letter ISO 3166-1 code for the country or territory |

AIR_PASSENGER_STATISTICS DATASET: It gives the number of passengers for different domestic and international airlines. It has a total of 11 attributes and 22870 records.

| Attributes | Data Types | Description |
|---|---|---|
| Passenger_Statistics_ID | Int | Primary Key(PK) |
| Activity Period | Int | The year and month at which passenger activity took place |

| | | |
|---|---|---|
| Operating Airline | Char | Airline name for the operator of aircraft with passenger activity |
| Operating Airline IATA Code | Char | The International Air Transport Association (IATA) two-letter designation for the Operating Airline. |
| Published Airline | Char | Airline name that issues the ticket and books revenue for passenger activity |
| Published Airline IATA Code | Char | The International Air Transport Association (IATA) two-letter designation for the Published Airline |
| GEO Summary | Char | Designates whether the passenger activity in relation to SFO arrived from or departed to a |
| | | location within the United States ("domestic"), or outside the United States ("international")without stops |
| GEO Region | Char | Provides a more detailed breakdown of the GEO Summary field to designate the region in the world where activity in relation to SFO arrived from or departed to without stops |
| Activity Type Code | Char | A description of the physical action a passenger took in relation to a flight, which includes boarding a flight ("enplanements"), getting off a flight ("deplanements") and transiting to another location ("in-transit") |

| | | |
|---|---|---|
| Price Category Code | Char | A categorization of whether a Published Airline is a low-cost carrier or not a low-cost carrier |
| Terminal | Char | The airport terminal designations at SFO where passenger activity took place |
| Boarding Area | Char | The airport boarding area designations at SFO where passenger activity took place |
| Passenger Count | Int | The number of monthly passengers associated with the above attribute fields. |

INTERNATIONAL _REPORT_DEPARTURE: This dataset provides the information on a specific airline for a pair of airports along with the number of flights from those airports. Three main columns record that record the number of flights is: Scheduled, Charter, and Total. It has a total of 16 attributes and approx. 68000 records.

| Attributes | Data Types | Description |
|---|---|---|
| Departure_ID | Int | Primary Key (PK) |
| date_dte | Date | Departure Date |
| Year | Date | Departure Year |
| Month | Date | Departure Month |
| Dept_usg_apt_id | Int | US Gateway airport ID |
| Dept_usg_apt | Char | US gateway IATA code |
| Dept_usg_wac | Int | US gateway world area code |
| Dept_fg_apt_id | Int | Foreign gateway airport ID |
| Dept_fg_apt | Char | Foreign gateway IATA code |
| Dept_fg_wac | int | Foreign gateway world area code |
| Dept_airlineid | Int | Airline ID |

| | | |
|---|---|---|
| Carrier | Char | Airline IATA code |
| Type | Char | Type of the metric |
| Dept_scheduled | Int | Metric flown by scheduled service operations |
| Dept_Charter | Int | Metric flown by charter operations |
| Dept_Total | Int | Total Metric flown by scheduled service and charter operations |

INTERNATIONAL_REPORT_PASSENGER: This dataset gives the information on the total number of passengers for each month and year between a pair of airports, as serviced by a particular airline.

| Attributes | Data Types | Description |
|---|---|---|
| Passenger_ID | Int | Primary Key (PK) |
| date_dte | Date | Date |
| Year | Date | Year |
| Month | Date | Month |
| usg_apt_id | Int | US Gateway airport ID |
| usg_apt | Char | US gateway IATA code |
| usg_wac | Int | US gateway world area code |
| fg_apt_id | Int | Foreign gateway airport ID |
| fg_apt | Char | Foreign gateway IATA code |
| fg_wac | Int | Foreign gateway world area code |
| airlineid | Int | Airline ID |
| Carrier | Char | Airline IATA code |
| Type | Char | Type of the metric |
| Scheduled | Int | Passengers flown by scheduled service operations |

| | | |
|---|---|---|
| Charter | Int | Passenger flown by charter operations |
| Total | Int | Total Passengers flown by scheduled service and charter operations |

# 3.0 Selection of Facts and Dimension Tables

**Facts tables**:

The Facts tables considered are as follows:

Airport (FactAirport)

Airline (FactAirlines)

**Dimensions tables**:

| Airport | Airlines |
|---|---|
| Countries (Dimcountries) | Routes (DimRoutes) |
| Passenger (Dimpassenger) | AirCarrier (DimAirCarrier) |
| Airport (DimAirport) | PassengerStats (DimPassengerStats) |
| International_ Departures(DimDepartures) | |
| International_Passengers(DimDPassengers) | |

**Types of measure**:

Airport

Fully Additive: Passenger Count

Airlines

Fully Additive: Passenger Count

# 4.0 Data Model

DimDepartures

| Department_ID (PK) |
| --- |
| Dept_date_dte |
| Dept_Year |
| Dept_Month |
| Dept_usg_apt |
| Dept_fg_apt |
| Carrier |
| Type |

DimPassengers

| Passenger_ID (PK) |
| --- |
| date_dte |
| Year |
| Month |
| usg_apt |
| fg_apt |
| Carrier |
| Type |

| Airport_ID (FK) |
| --- |
| Country_ID (FK) |
| Passenger_ID (FK) |
| Department_ID (FK) |
| Latitude |
| Longitude |
| Altitude |
| Timezone |
| Passenger_count |
| Dept_usg_apt_id |
| Dept_usg_wac |
| Dept_fg_apt_id |
| Dept_fg_wac |
| Dept_airlineid |
| Dept_carriergroup |
| Dept_scheduled |
| Dept_charter |
| Dept_total |
| usg_apt_id |
| usg_wac |
| fg_apt_id |
| fg_wac |
| airlineid |
| carriergroup |
| scheduled |
| charter |
| total |

DimAirport

| Airport_ID (PK) |
| --- |
| Name |
| City |
| Country |
| IATA |
| ICAO |
| DST |
| Tz database timezone |
| Type |
| Source |

FactAirport

DimCountries

| Country_ID (PK) |
| --- |
| name |
| iso_code |
| dafif_code |

DimRoutes

| Routes_ID (PK) |
| --- |
| Airline_ID(FK) |
| Airline |
| Source Airport |
| Source Airport_ID |
| Destination Airport |
| Destination Airport_ID |
| Codeshare |
| Equipments |

FactAirlines

| Airline_ID (FK) |
| --- |
| Route_ID (FK) |
| Passenger_statistics_ID (FK) |
| Source Airport_ID |
| Destination Airport_ID |
| Stops |
| Passenger Count |

DimAirCarrier

| Airline_ID (PK) |
| --- |
| Name |
| Alias |
| IATA |
| ICAO |
| Callsign |
| Country |
| Active |

| Passenger_Statistics_ID (PK) |
| --- |
| Activity Year |
| Activity Month |
| Operating Airline |
| Operating Airline IATA Code |
| Published Aiirline |
| Published Airline IATA Code |
| GEO Summary |
| GEO Region |
| Activity Type Code |
| Price Category Code |
| Terminal |
| Boarding Area |

DimPassengerStats

## 5.0 Data Flow:

We would Consider Flat File Source Connection, the source would be in CSV Format (Aviation Datasets from Kaggle and Flight.org.) In the staging area, we would have a condition block to load bulk data. After a bulk load, we will have the condition to load weekly data, till the rest of the data is transferred into our OLE Db Destination. In destination, we will create facts and dimension tables as mentioned below.

### 5.0.1 Basic Data Flow



### 5.0.2 High-level View

## 6.0 Loading and staging

As discussed in the previous sections, we will initially load the data from the source datasets, through the appropriate source tables.

We load the data initially into the Staging table.

We process the data by using tasks present in AWS Glue Databrew like Data Conversion to check and convert the data types appropriately, adding Slowly Changing Dimensions to add the historical values to attributes like Callsign, and PassengerCount.

With the changes reflected into the appropriate fields along with the StartDate, End Date and the Flag, we load the data into the Star/ Snowflake Schema of the destination tables.

We consider look up task to probe the IATA and examine the modifications of data through Slowly Changing Dimension tasks.

S3 Bucket created for storing the CSV file: We created a two S3 buckets for storing the data.

The 1st bucket was created to store the raw data and the second bucket was created to store the cleaned the data.

1st S3 bucket: dwbi-project04-stage

2nd S3 bucket: dwbi-project04-new

After the raw data sets were stored into the S3 bucket, we cleaned the datasets to remove all the unnecessary data using AWS Glue Databrew.

In the above fig you can see the projects created for all the data sets. For each Dataset we were required to create separate projects and then we ran recopied on each of the projects.



After that we removed unnecessary columns for all the datasets and recipes were created for each datasets.

For each individual project different recipes were created. Recipes are nothing but the data cleaning process that were done to obtain clean data set. After each recipe a JOB was run. The JOB implemented the recipe to the whole dataset and then it was store to the new S3 bucket i.e., dwbi-project04-new.

Similarly, Jobs were run for all the recipes to store the cleaned data in the new bucket.



For each recipe different JOB was created and it was run to obtain all the clean datasets. The fig above shows 7 JOB's that were created for all the data sets.

To check if the JOBs were successful, we pulled the clean data to Athena and we did a query to check if the data is being populated and if all the unwanted data are removed from the datasets.

In order to load the datasets to the data warehouse we used the ETL tool data Glue. In Data glue we created different crawlers to add the tables for all the datasets. As we can see from the image below, 14 different crawlers were created. In those 14 crawlers, 7 crawlers were created to pull the data form the source table i.e., the S3 bucket and the other 7 crawlers were created in to connect it to the destination i.e Redshift.

The above fig shows the properties of the Redshift data crawler that was created to connect the Air-Pass-Stats-Redshift table in the redshift to the source. For this we used the JDBC connection and provided the proper service role and included the path for the connection. The image below shows the connection to the Redshift.
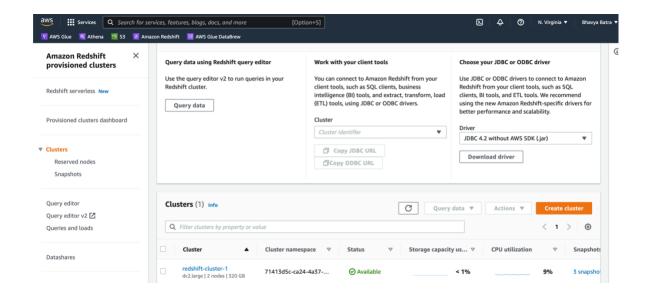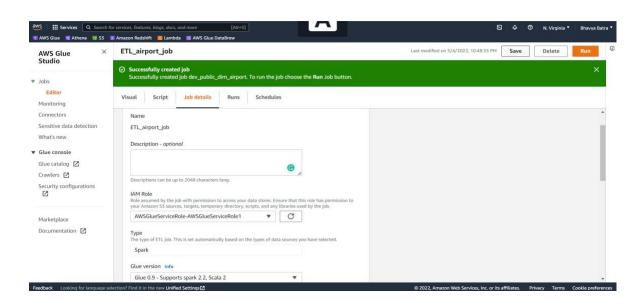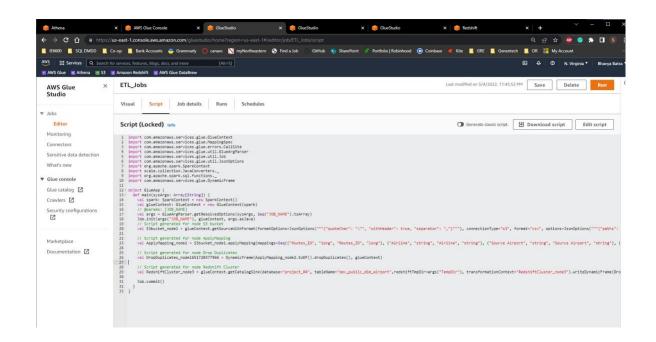


# 7.0 Creation of Data Warehouse

After the crawlers were created, we were needed to create jobs. The Jobs were used in order to map the data table form the source to the data warehouse. The diagram below shows how the normal schema was created for a certain job. S3 bucket was connected to the redshift via mapping and drop duplicates transformation tools.
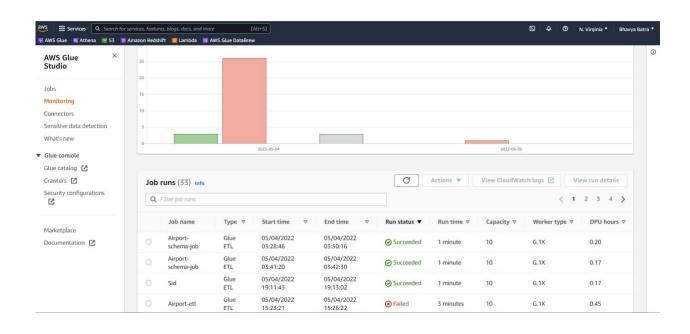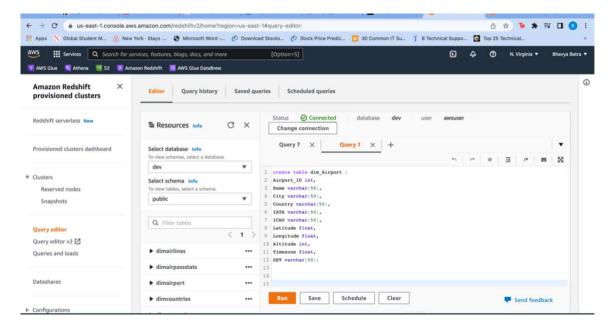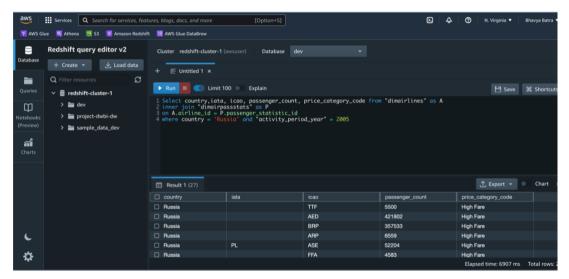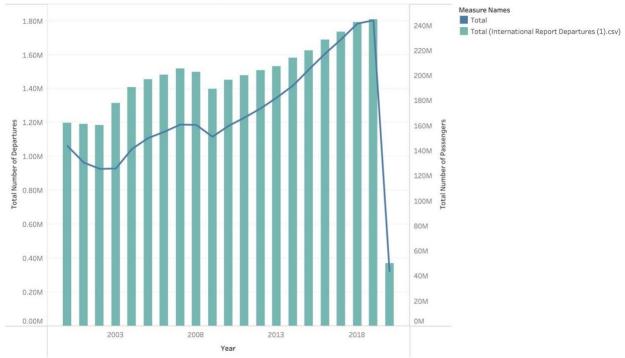
Ran the queries in the query editor to check if the data has been loaded into the new tables.

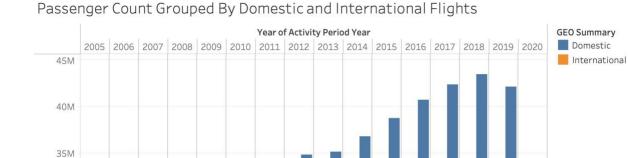# 8.0 Visualization Expected:

Correlation Between Number of Passenger & Airlines Over The Years



The trends of Total (International Report Departures (1).csv) and Total for Year (International Report Departures (1).csv) Year.
Color shows details about Total (International Report Departures (1).csv) and Total. The data is filtered on Year (International Report Departures (1).csv), which ranges from 01/01/2000 to 01/01/2020.

- This visualization shows the relationship between the number of passengers and the number of flights departed from the year 2000 to 2020.
- It shows a strong positive correlation between the two attributes analyzed here.
- There is a sharp decrease in the number of departures and passengers in the year 2020 that can be speculated because of covid-19.

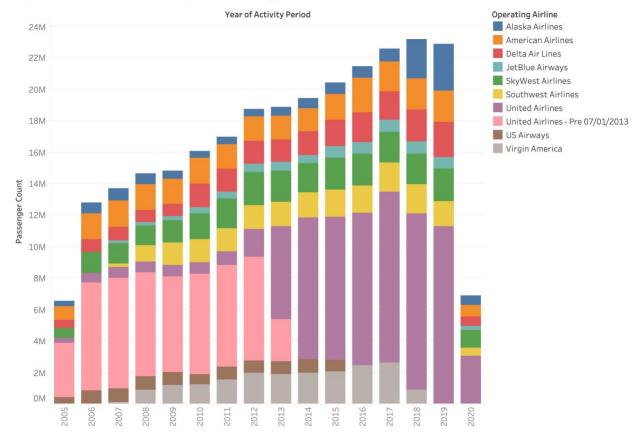## Passenger Count Grouped By Domestic and International Flights



Sum of Passenger Count for each GEO Summary broken down by Activity Period Year Year. Color shows details about GEO Summary.

- This visualization shows the comparison between the number of passengers in domestic flights and international flights departing from the United States between the years 2005 to 2020.
- The number of passengers in international flights is more than twice the number of passengers in domestic flights every year. .
- Each year, the number of passengers increases for both domestic flights and international flights.
- We can conclude that the economic growth in a country influences the growth in air travel

## Top 10 Airlines With The Highest Passenger Count From 2015-2020

**Year of Activity Period**



Sum of Passenger Count for each Activity Period Year. Color shows details about Operating Airline. The view is filtered on Operating Airline, which keeps 10 of 95 members.

- This visualization shows the top 10 highest performing airports in USA based on the passenger count between 2005 and 2020.
- As we can see, the visualization shows that United Airlines has monopolized the industry with almost half of all the total passengers choosing United Airlines as their preferred airlines.
- The other airlines have been performing consistently since the past two decades, except for Virgin America having a sudden increase and decrease between 2008 and 2018.

## Countries With The Largest Number Of Airports



Count of IATA_Code
334          1,512

© 2022 Mapbox © OpenStreetMap

Map based on Longitude (generated) and Latitude (generated). Color shows count of IATA_Code. Details are shown for Country. The view is filtered on count of IATA_Code, which includes values greater than or equal to 302.

- This visualization shows the highest number of airports in the world based on our sample of airports dataset.
- It counts the number of unique IATA codes and shows which country has the highest number of airports(between 300 and 1512).
- As we can see, the highest number of airports are located in ISA, Canada and Australia. We can predict that this has to do with the country's economy, size, tourism, etc. Although a lot of factors influence the number of airports in a country, airports are a symbol of growth and prosperity of the country.
- This visualization does not state the absolute facts, but only represents the major ROI of our dataset.

## 9.0 Modification History:

| Version | Date | Author | Description of Edits |
| --- | --- | --- | --- |
| 1 | 03/23 | Team | Project Proposal (Draft) |
| 2 | 04/06 | Team | Second Draft |
| 2.1 | 03/31 | Bhavya | Data source |
| 2.2 | 04/01 | Ronen | Facts and Dimension |
| 2.3 | 04/02 | Aishwarya | ER Diagram |
| 2.4 | 04/03 | Siddhesh | Word Document |
| 3 | 04/19 | Team | Final Draft |
| 3.1 | 04/19 | Ronen | Data Set Description and Cleaning, Word Document |
| 3.2 | 04/20 | Siddhesh | SCD Types, Word Document |
| 3.3 | 04/21 | Aishwarya | ER Diagram updated, Data set Description |
| 3.4 | 04/21 | Bhavya | Data set Description, Word Document |
| 4 | 04/28 | Team | Loading Data to Data Warehouse |
| 4.1 | 04/28 | Bhavya | Loading data into AWS S3 and cleaning the data using AWS Databrew |
| 4.2 | 04/29 | Ronen | Loading the cleaned data into S3 |
| 4.3 | 04/30 | Siddhesh | Loading Data from S3 to Redshift |
| 4.4 | 05/02 | Aishwarya | Creating Visualization |
| 5 | 05/04 | Team | Creating the Report and recording the Video |