# ETL Project

## By Anviksha Singh

anviksha.singh0110s@gmail.com

## Step 3: Redshift Set Up and Data Loading

# Creation of a Redshift Cluster

## Screenshots of the configuration of the Redshift cluster that you have created:

## Screenshot of data loaded in S3

# Setting up a database in the Redshift cluster and running queries to create the dimension and fact tables



## Queries to create the various dimension and fact tables with appropriate primary and foreign keys:

Step 1 : Create Schema

```
create schema atm_trans_data;
```

Step 2 : Create dimension tables

DIM_LOCATION

```
create table atm_trans_data.DIM_LOCATION
(
 location_id int not null DISTKEY SORTKEY,
location varchar(50),
streetname varchar(255),
street_number int,
zipcode int,
lat decimal(10,3),
lon decimal(10,3),
PRIMARY KEY(location_id)
);
```

**DIM_CARD TYPE**

```
create table atm_trans_data.DIM_CARD_TYPE
(
card_type_id int not null DISTKEY SORTKEY,
card_type varchar(30),
PRIMARY KEY(card_type_id)
);
```

**DIM_ATM**

```
create table atm_trans_data.DIM_ATM
(
atm_id int not null DISTKEY SORTKEY,
atm_number varchar(20),
atm_manufacturer varchar(50),
atm_location_id int,
PRIMARY KEY(atm_id),
FOREIGN KEY(atm_location_id) references atm_trans_data.DIM_LOCATION(location_id)
);
```

**DIM_DATE**

```
create table atm_trans_data.DIM_DATE
(
date_id int not null DISTKEY SORTKEY,
full_date_time timestamp,
year int,
month varchar(20),
day int,
hour int,
weekday varchar(20),
PRIMARY KEY(date_id)
);
```

Step 3 : Create Fact table

```sql
create table atm_trans_data.FACT_ATM_TRANS
(
trans_id bigint not null DISTKEY SORTKEY,
atm_id int,
weather_loc_id int,
date_id int,
card_type_id int,
atm_status varchar(20),
currency varchar(10),
service varchar(20),
transaction_amount int,
message_code varchar(225),
message_text varchar(225),
rain_3h decimal(10,3),
clouds_all int,
weather_id int,
weather_main varchar(50),
weather_description varchar(255),
PRIMARY KEY(trans_id),
FOREIGN KEY(weather_loc_id) references atm_trans_data.DIM_LOCATION(location_id),
FOREIGN KEY(atm_id) references atm_trans_data.DIM_ATM(atm_id),
FOREIGN KEY(date_id) references atm_trans_data.DIM_DATE(date_id),
FOREIGN KEY(card_type_id) references atm_trans_data.DIM_CARD_TYPE(card_type_id)
);
```

# Loading data into a Redshift cluster from Amazon S3 bucket

## Queries to copy the dimension dim_atm from S3 buckets to the Redshift cluster

A. <u>S3 Bucket – Dim_atm</u>



B. <u>Query to copy the dim_atm from S3 buckets to the Redshift cluster</u>

```
COPY "upgrad-anvi".atm_trans_data.dim_atm FROM
's3://etlfactdimbuck/dim_atm/part-00000-d72e9453-a7dc-41a2-97a0-
0eb5e6ffb405-c000.csv' IAM_ROLE
'arn:aws:iam::128629367978:role/myRedshiftRole' FORMAT AS CSV DELIMITER
',' QUOTE '"' REGION AS 'us-east-1'
```

C. <u>Loaded Data</u>

```
68    Select * from atm_trans_data.dim_atm;
```

**Result 1 (100)**

| atm_id | atm_number | atm_manufacturer | atm_location_id |
|--------|-----------|------------------|-----------------|
| 0 | 40 | Diebold Nixdorf | 86 |
| 3 | 28 | NCR | 33 |
| 8 | 81 | NCR | 81 |
| 12 | 45 | NCR | 68 |
| 13 | 79 | NCR | 42 |
| 17 | 86 | NCR | 104 |
| 21 | 12 | NCR | 69 |

# Queries to copy the dimension dim_card_type from S3 buckets to the Redshift cluster

A. <u>S3 Bucket – Dim_card_type</u>



B. <u>Query to copy the dim_card_type from S3 buckets to the Redshift cluster</u>

```
COPY "upgrad-anvi".atm_trans_data.dim_card_type FROM
's3://etlfactdimbuck/dim_card_type/part-00000-3e9ed684-85f8-430b-b1af-
5dce4951ec47-c000.csv' IAM_ROLE
'arn:aws:iam::128629367978:role/myRedshiftRole' FORMAT AS CSV DELIMITER
',' QUOTE '"' REGION AS 'us-east-1'
```

C. <u>Loaded Data</u>

```
69    Select * from atm_trans_data.dim_card_type;
```

Result 1 (12)

| card_type_id | card_type |
|---|---|
| 1 | Mastercard - on-us |
| 6 | HÃƒÂ¡vekort - on-us |
| 11 | Dankort |
| 5 | Visa Dankort |
| 2 | HÃƒÂ¡vekort |
| 4 | Dankort - on-us |
| 7 | CIRRUS |
| 9 | Maestro |
| 10 | MasterCard |
| 0 | Visa Dankort - on-us |
| 3 | VisaPlus |
| 8 | VISA |

# Queries to copy the dimension dim_date from S3 buckets to the Redshift cluster

A. <u>S3 Bucket – Dim_date</u>



B. <u>Query to copy the dim_date from S3 buckets to the Redshift cluster</u>

```
COPY "upgrad-anvi".atm_trans_data.dim_date FROM
's3://etlfactdimbuck/dim_date/part-00000-73506087-a16f-4d00-b385-
5220080c8029-c000.csv' IAM_ROLE
'arn:aws:iam::128629367978:role/myRedshiftRole' FORMAT AS CSV DELIMITER
',' QUOTE '"' TIMEFORMAT 'auto' REGION AS 'us-east-1'
```
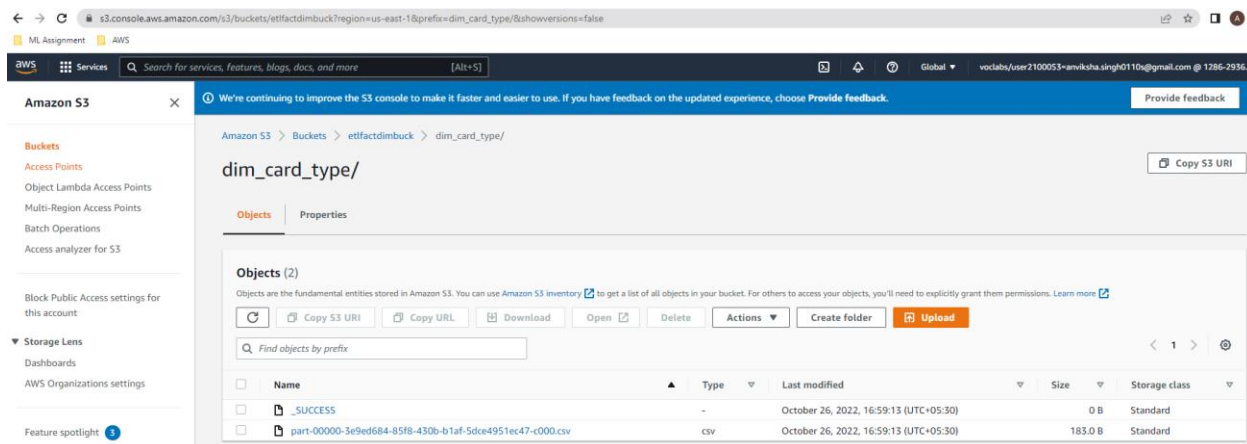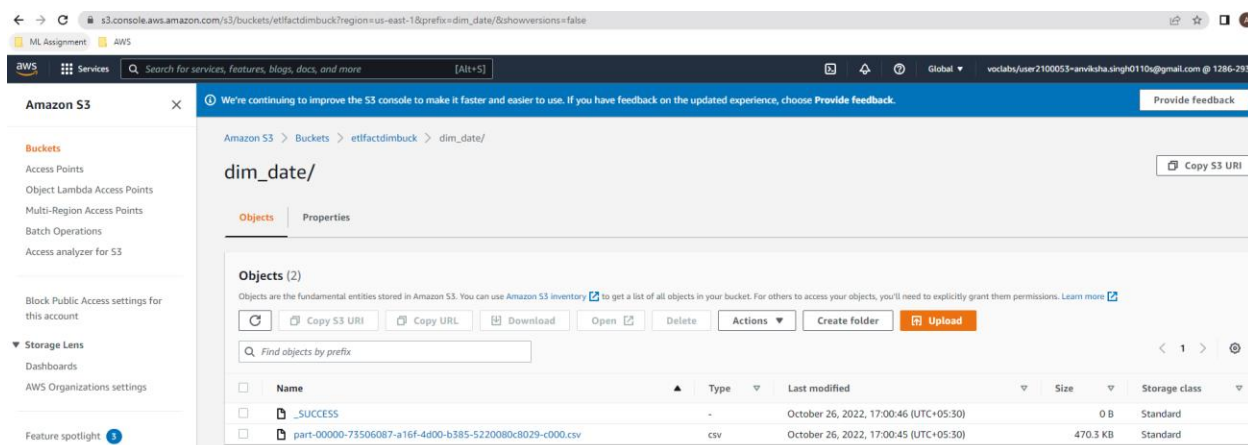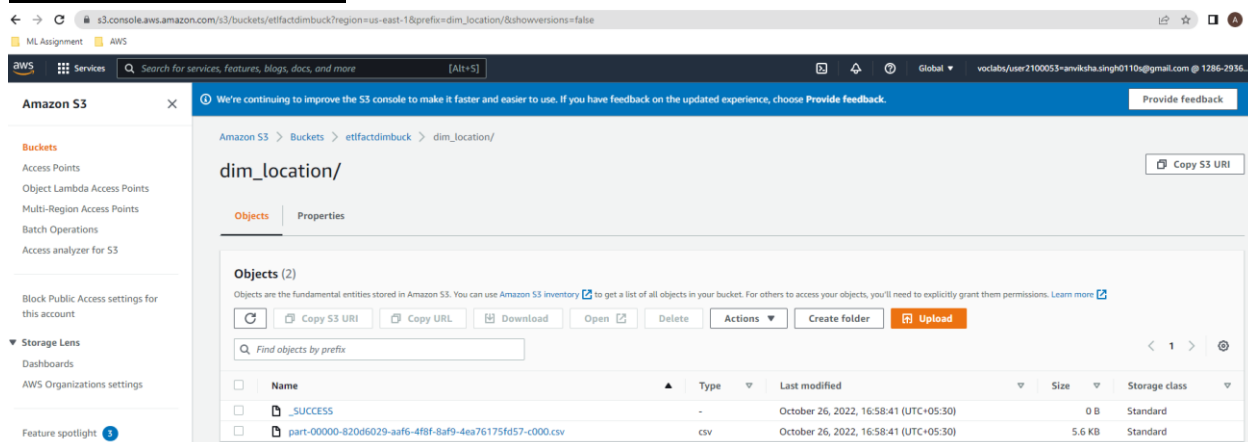
C. <u>Loaded Data</u>

```
70    Select * from atm_trans_data.dim_date;
```

Result 1 (100)

| date_id | full_date_time | year | month | day | hour | weekday |
|---------|----------------|------|-------|-----|------|---------|
| 5 | 2017-01-26 20:00:00 | 2017 | January | 26 | 20 | Thursday |
| 15 | 2017-02-04 23:00:00 | 2017 | February | 4 | 23 | Saturday |
| 22 | 2017-03-26 19:00:00 | 2017 | March | 26 | 19 | Sunday |
| 23 | 2017-01-20 12:00:00 | 2017 | January | 20 | 12 | Friday |
| 24 | 2017-03-28 02:00:00 | 2017 | March | 28 | 2 | Tuesday |
| 27 | 2017-01-13 23:00:00 | 2017 | January | 13 | 23 | Friday |
| 28 | 2017-02-12 06:00:00 | 2017 | February | 12 | 6 | Sunday |
| 33 | 2017-02-07 01:00:00 | 2017 | February | 7 | 1 | Tuesday |
| 35 | 2017-01-10 20:00:00 | 2017 | January | 10 | 20 | Tuesday |
| 39 | 2017-01-18 20:00:00 | 2017 | January | 18 | 20 | Wednesday |
| 42 | 2017-02-09 23:00:00 | 2017 | February | 9 | 23 | Thursday |
| 44 | 2017-03-31 07:00:00 | 2017 | March | 31 | 7 | Friday |
| 46 | 2017-01-06 00:00:00 | 2017 | January | 6 | 0 | Friday |

# Queries to copy the dimension dim_location from S3 buckets to the Redshift cluster

A. <u>S3 Bucket – Dim_location</u>



B. <u>Query to copy the dim_location from S3 buckets to the Redshift cluster</u>

```
COPY "upgrad-anvi".atm_trans_data.dim_location FROM
's3://etlfactdimbuck/dim_location/part-00000-820d6029-aaf6-4f8f-8af9-
4ea76175fd57-c000.csv' IAM_ROLE
'arn:aws:iam::128629367978:role/myRedshiftRole' FORMAT AS CSV DELIMITER
',' QUOTE '"' REGION AS 'us-east-1'
```
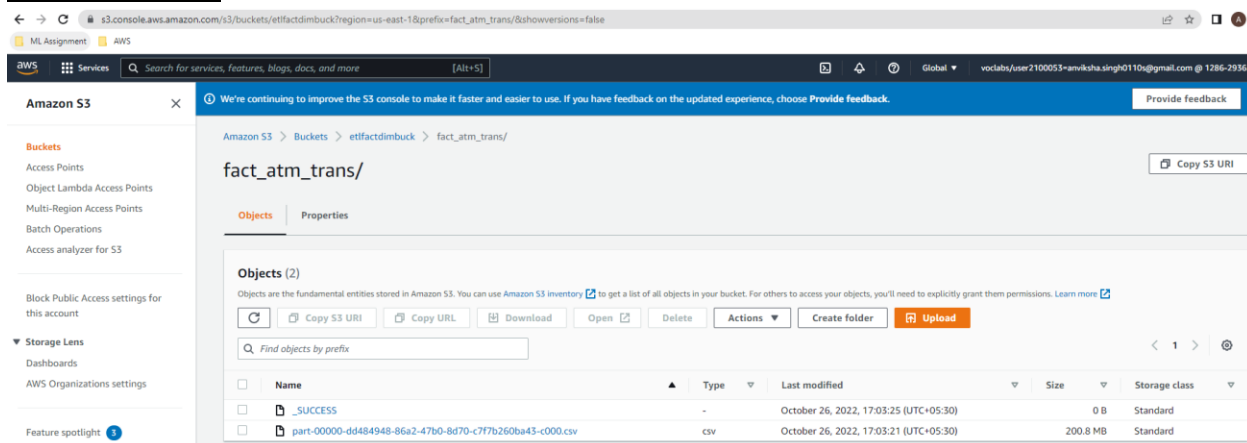
C. <u>Loaded Data</u>

```
71   Select * from atm_trans_data.dim_location;
```

**Result 1 (100)**

| location_id | location | streetname | street_number | zipcode | lat | lon |
|---|---|---|---|---|---|---|
| 1 | Skelagervej 15 | Skelagervej | 15 | 9000 | 57.023 | 9.891 |
| 6 | Esbjerg | Strandbygade | 20 | 6700 | 55.468 | 8.44 |
| 11 | NÃƒÂ¦stved | Farimagsvej | 8 | 4700 | 55.233 | 11.763 |
| 14 | Slagelse | Mariendals ... | 29 | 4200 | 55.398 | 11.342 |
| 20 | Brugsen ANS | SÃƒÂ¸nder... | 14 | 8643 | 56.306 | 9.594 |
| 32 | Sauersvej | Fridtjof Nan... | 2 | 9210 | 57.023 | 9.94 |
| 34 | Aalborg Storce... | Hobrovej | 452 | 9200 | 57.005 | 9.876 |
| 36 | Bispensgade | Bispensgade | 35 | 9800 | 57.453 | 9.996 |
| 48 | Frederiksberg | Gammel Ko... | 157 | 1850 | 55.677 | 12.537 |
| 49 | Storcenter indg... | Hobrovej | 452 | 9200 | 57.005 | 9.876 |
| 50 | HjÃƒÂ¸rring | ÃƒÂ¸esterg... | 8 | 9800 | 57.459 | 9.988 |
| 51 | NÃƒÂ¦stved | Farimagsgade | 8 | 4700 | 55.69 | 12.575 |
| 53 | Skallerup Klit | Nordre Klitvej | 21 | 9800 | 57.494 | 9.838 |
| 55 | ÃƒÂ¸esterÃƒÂ... | ÃƒÂ¸esterÃƒ... | 12 | 9000 | 57.049 | 9.922 |
| 58 | Intern Frederik | Danmarksg | 48 | 9900 | 57.441 | 10.537 |

# Queries to copy the fact from S3 buckets to the Redshift cluster

A. S3 Bucket – fact



B. Query to copy the fact from S3 buckets to the Redshift cluster

```
COPY "upgrad-anvi".atm_trans_data.fact_atm_trans FROM
's3://etlfactdimbuck/fact_atm_trans/part-00000-dd484948-86a2-47b0-8d70-
c7f7b260ba43-c000.csv' IAM_ROLE
'arn:aws:iam::128629367978:role/myRedshiftRole' FORMAT AS CSV DELIMITER
',' QUOTE '"' REGION AS 'us-east-1'
```

C. Loaded Data

**Performance after all the query execution:**