

🔑 main ▾

...

Data_Lake_Project / README.md



Abhaycl Update README.md

🕒 History

👤 1 contributor

Raw

Blame



331 lines (202 sloc) | 16.4 KB

Data Lake Project Starter Code

The objective of this project is to apply what we've learned on Spark and data lakes to build an ETL pipeline for a data lake hosted on S3.

How to run the program with your own code

For the execution of your own code, we head to the Project Workspace.

Running the etl.py script in a terminal to process the all the datasets. The script connects to the Sparkify database, extracts and processes the log_data and song_data, and loads data into the five tables.

```
python etl.py
```

The summary of the files and folders within repo is provided in the table below:

File/Folder	Definition
-------------	------------

File/Folder	Definition
data/*	Folder that contains all the json files with the data used in this project.
images/*	Folder containing the images of the project.
dl.cfg	Contains our AWS credentials.
etl.py	Contains all the processing of reads data from S3, processes that data using Spark, and writes them back to S3.
README.md	Contains the project documentation.
README.pdf	Contains the project documentation in PDF format.

Steps to complete the project:

1. To complete the project, you will need to load data from S3, process the data into analytics tables using Spark, and load them back into S3. You'll deploy this Spark process on a cluster using AWS..

Rubric Points

Here I will consider the rubric points individually and describe how I addressed each point in my implementation.

Scenario.

A music streaming startup, Sparkify, has grown their user base and song database even more and want to move their data warehouse to a data lake. Their data resides in S3, in a directory of JSON logs on user activity on the app, as well as a directory with JSON metadata on the songs in their app.

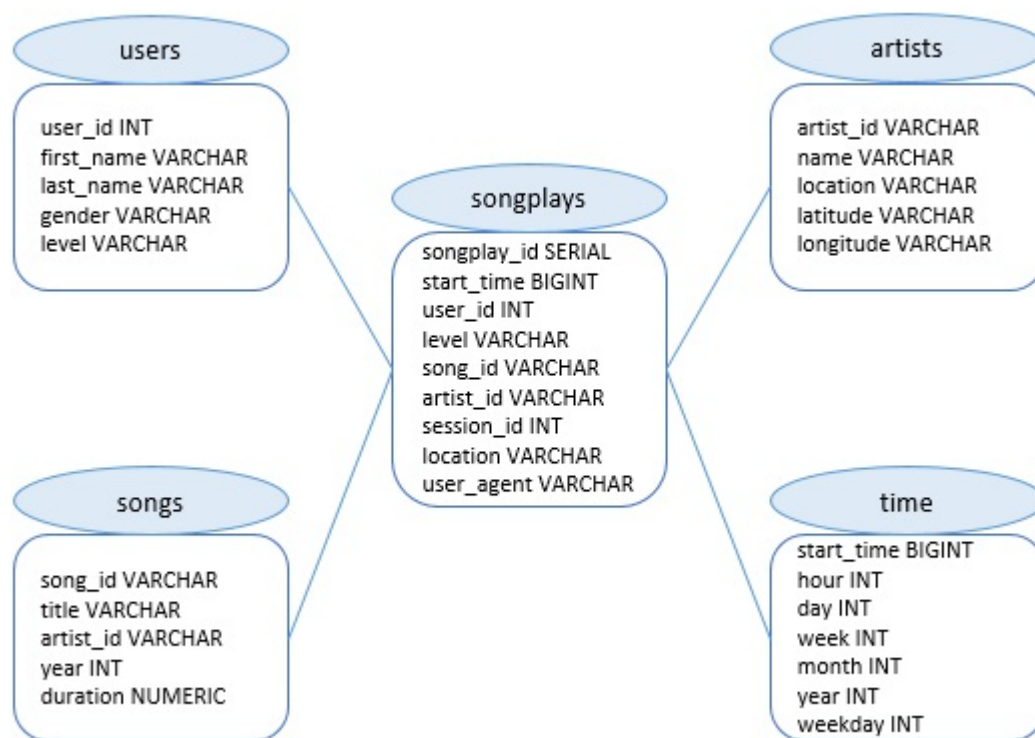
As their data engineer, you are tasked with building an ETL pipeline that extracts their data from S3, processes them using Spark, and loads the data back into S3 as a set of dimensional tables. This will allow their analytics team to continue finding insights in what songs their users are listening to.

You'll be able to test your database and ETL pipeline by running queries given to you by the analytics team from Sparkify and compare your results with their expected results.

Database Schema for Sparkify.

In our scenario, we have one dataset of songs where each file is in JSON format and contains metadata about a song and the artist of that song. The second dataset consists of log files in JSON format and simulates activity logs from a music streaming app (based on specified configurations).

Using the song and log datasets, the database schema looks like this:



Star Schema.

A star schema is the simplest style of data mart schema. The star schema consists of one or more fact tables referencing to any number of dimension tables. It has some advantages like fast aggregation for analytics, simple queries for JOINS, etc.

ETL Pipeline

Extract, transform, load (ETL) is the general procedure of copying data from one or more sources into a destination system which represents the data differently from, or in a different context than, the sources.

Project Datasets

We'll be working with two datasets that reside in S3. Here are the S3 links for each:

- Song data: `s3://udacity-dend/song_data`
- Log data: `s3://udacity-dend/log_data`

Song Dataset

The first dataset is a subset of real data from the [Million Song Dataset](#). Each file is in JSON format and contains metadata about a song and the artist of that song. The files are partitioned by the first three letters of each song's track ID. For example, here are filepaths to three files in this dataset.

```
song_data/A/A/A/TRAAAW128F429D538.json
song_data/A/A/B/TRAABCL128F4286650.json
song_data/A/A/C/TRAACCG128F92E8A55.json
```

And below is an example of what a single song file, `TRAAAW128F429D538.json`, looks like.

```
{"num_songs": 1, "artist_id": "ARD7TVE1187B99BFB1", "artist_latitude": null, "artist
```

Log Dataset

The second dataset consists of log files in JSON format generated by this event simulator based on the songs in the dataset above. These simulate app activity logs from an imaginary music streaming app based on configuration settings.

The log files in the dataset we'll be working with are partitioned by year and month. For example, here are filepaths to three files in this dataset.

```
log_data/2018/11/2018-11-01-events.json
log_data/2018/11/2018-11-02-events.json
log_data/2018/11/2018-11-03-events.json
```

And below is an example of what the data in a log file, `2018-11-01-events.json`, looks like.

```
{"artist": null, "auth": "Logged In", "firstName": "Walter", "gender": "M", "itemInSession":
```

Prerequisites

It is necessary that we have an AWS account. The following tasks must be realized within AWS:

1. In the IAM area you have to create a user with read and write permissions or full access to the S3 buckets area.
2. In the S3 area you have to create an S3 bucket, in our case we call it udacity-datalake-project-abhaycl.

The image shows two screenshots of the AWS S3 console. The top screenshot is the 'Create bucket' page. The bottom screenshot is the 'Buckets' list page.

Top Screenshot: Create bucket

The 'Create bucket' page shows the following configuration:

- Bucket name:** udacity-datalake-project-abhaycl
- AWS Region:** US West (Oregon) us-west-2
- Copy settings from existing bucket - optional:** Choose bucket
- Block Public Access settings for bucket:** ☐ Block all public access

Bottom Screenshot: Buckets (1)

The 'Buckets (1)' page shows a table with one bucket:

Name	AWS Region	Access	Creation date
udacity-datalake-project-abhaycl	US West (Oregon) us-west-2	Objects can be public	February 17, 2021, 06:05:15 (UTC+01:00)

Possible help commands from the console:

- Empty the content of the S3 bucket: `aws s3 rm s3://udacity-datalake-project-abhaycl --recursive`
- Delete the previously created S3 bucket: `aws s3 rb s3://udacity-datalake-project-abhaycl --force`

Note: If we have an AWS Free Tier account, each month only allows 2000 Put, Copy, Post or List Requests of Amazon S3, this may limit the reading of the song data files from the source S3 bucket (s3a://udacity-dend/song_data/) on our destination S3 bucket (s3://udacity-datalake-project-abhaycl).

Process song data (song_data directory).

We will perform ETL on the files in song_data directory to create two dimensional tables: songs table and artists table.

This is what a songs file looks like:

```
{"num_songs": 1, "artist_id": "ARD7TVE1187B99BFB1", "artist_latitude": null, "artist
```

We will proceed to read the data from the source S3 bucket, in the image below we can check the result of the task.

```
df = spark.read.json(song_data)
```

```
--- Reading song data from json files. --- 2021-02-17
- Shows the scheme and some song data records...
root
 |-- artist_id: string (nullable = true)
 |-- artist_latitude: double (nullable = true)
 |-- artist_location: string (nullable = true)
 |-- artist_longitude: double (nullable = true)
 |-- artist_name: string (nullable = true)
 |-- duration: double (nullable = true)
 |-- num_songs: long (nullable = true)
 |-- song_id: string (nullable = true)
 |-- title: string (nullable = true)
 |-- year: long (nullable = true)

+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| artist_id|artist_latitude| artist_location|artist_longitude| artist_name| duration|num_songs| song_id| title|year|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|ARTC11V1187B9A4858| 51.4536|Goldsmith's Colle...| -0.01802| The Bonzo Dog Band|301.40036| 1|SOAFBCP12A8C13C7D|King Of Scurf (20...|1972|
|ARA23X01187B9AF18F| 40.57885|Carteret, New Jersey| -74.21956| The Smithereens| 192.522| 1|SOXTJDS12AF72A2SE5|Drown In My Own T...| 0|
|ARSVTNL1187B9A2A91| 51.50632| London, England| -0.12714| Jonathan King|129.85424| 1|SOEKAZG12A8018837E|I'll Slap Your Fa...|2001|
|AR73AI01187B9AD57B| 37.77916| San Francisco, CA| -122.42005| Western Addiction|118.07302| 1|SOQFWCR12A6D4FB2A3|A Poor Recipe For...|2005|
|ARXQBR11187B98A2CC| null| Liverpool, England| null|Frankie Goes To H...|821.05424| 1|SOBRKGM12A8C139TF6|Welcome to the Pl...|1985|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows
```

For the songs table, we will extract and transform the song data using only the columns corresponding to the songs table suggested in the project specifications, sorting the records by song ID and without duplicates.

```
songs_table = df.select("song_id", "title", "artist_id", "year", "duration") \
    .orderBy("song_id") \
    .drop_duplicates()
```

```

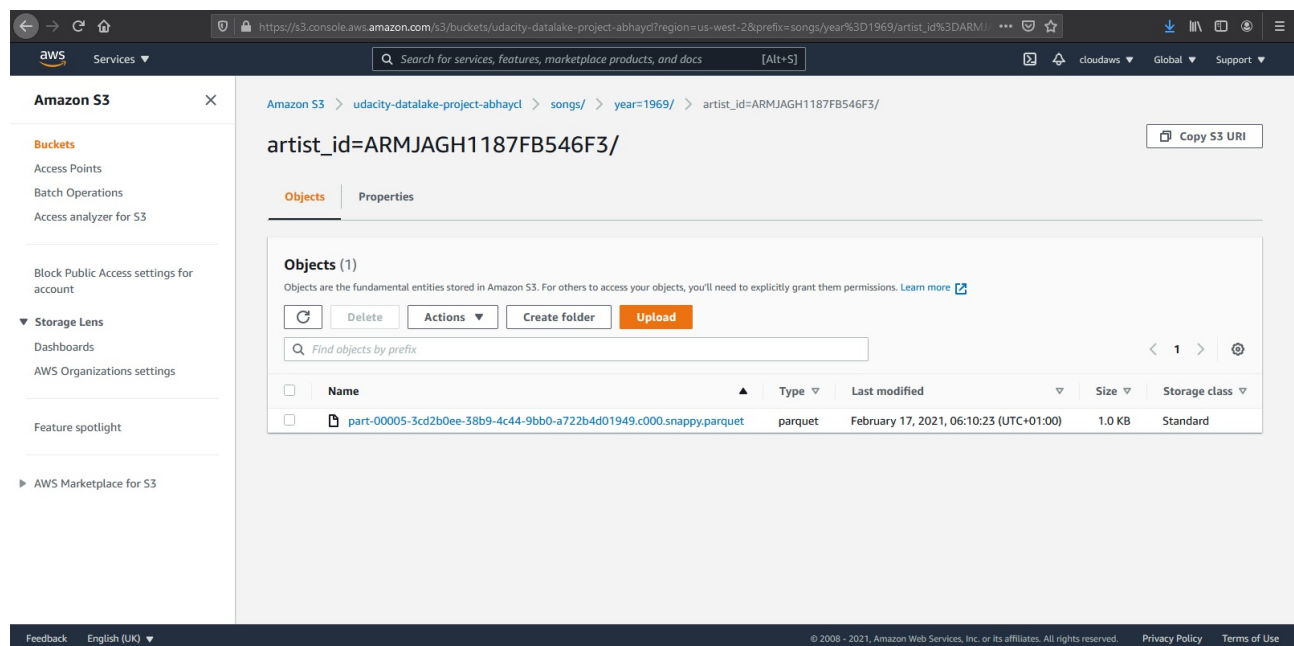
--- Extraction of columns to create songs table. ---
- Shows the scheme and some songs table records...
root
 |-- song_id: string (nullable = true)
 |-- title: string (nullable = true)
 |-- artist_id: string (nullable = true)
 |-- year: long (nullable = true)
 |-- duration: double (nullable = true)

+-----+-----+-----+-----+-----+
| song_id| title| artist_id| year| duration|
+-----+-----+-----+-----+
|SOABWAP12A8C13F82A| Take Time| AR5LMPY1187FB573FE| 1978| 258.89914|
|SOAFBCP12A8C13CC7D| King Of Scurf (20...| ARTC1LV1187B9A4858| 1972| 301.40036|
|SOAPERH12A58A787DC| The One And Only ...| ARZ5HOP1187B98A1DD| 0| 230.42567|
|SOBLFFE12AF72AA5BA| Scream| ARJNUY12298900C91| 2009| 213.9424|
|SOBRKGM12A8C139EF6| Welcome to the Pl...| ARXQBR1187B98A2CC| 1985| 821.05424|
+-----+-----+-----+-----+
only showing top 5 rows

```

We proceed to load our data into our destination S3 bucket to write the song table in parquet files partitioned by year and artist.

```
songs_table.write.partitionBy("year", "artist_id").parquet(os.path.join(output_d
```



For the artist table, we will extract and transform the song data using only the columns corresponding to the artist table suggested in the project specifications, sorting the records by artist ID and without duplicates.

```

artists_table = df.selectExpr("artist_id", "artist_name as name", "artist_locati
.orderBy("artist_id") \

```



```
.drop_duplicates()
```

```
--- Extraction of columns to create artists table. ---
- Shows the scheme and some artists table records...
root
|-- artist_id: string (nullable = true)
|-- name: string (nullable = true)
|-- location: string (nullable = true)
|-- latitude: double (nullable = true)
|-- longitude: double (nullable = true)

+-----+-----+-----+-----+-----+
|      artist_id|      name|      location|latitude| longitude|
+-----+-----+-----+-----+-----+
|AR0MWD61187B9B2B12|International Noi...|      |      null|      null|
|AR10USD1187B99F3F1|Tweeterfriendly M...|Burlington, Ontar...|      null|      null|
|AR1C2IX1187B99BF74|      Broken Spindles|      |      null|      null|
|AR1KTV21187B9ACD72|      Cristina|      California - LA|34.05349|-118.24532|
|AR5LMPY1187FB573FE|      Chaka Khan_ Rufus|      Chicago, IL|41.88415|-87.63241|
+-----+-----+-----+-----+-----+
only showing top 5 rows
```

We proceed to load our data into our destination S3 bucket to write the artist table in parquet files.

```
artists_table.write.parquet(os.path.join(output_data, "artists"), mode="overwrite")
```

Amazon S3 console view of the `artists/` directory. The `Objects` tab is active, showing 25 objects. The first object is `_SUCCESS` (0 B). The following four objects are parquet files:

Name	Type	Last modified	Size	Storage class
<code>_SUCCESS</code>	-	February 17, 2021, 06:15:26 (UTC+01:00)	0 B	Standard
<code>part-00000-ff82e716-c02d-4bba-b5ae-84964477feaa-c000.snappy.parquet</code>	parquet	February 17, 2021, 06:14:14 (UTC+01:00)	1.5 KB	Standard
<code>part-00001-ff82e716-c02d-4bba-b5ae-84964477feaa-c000.snappy.parquet</code>	parquet	February 17, 2021, 06:14:17 (UTC+01:00)	1.6 KB	Standard
<code>part-00002-ff82e716-c02d-4bba-b5ae-84964477feaa-c000.snappy.parquet</code>	parquet	February 17, 2021, 06:14:21 (UTC+01:00)	1.3 KB	Standard
<code>part-00003-ff82e716-c02d-4bba-b5ae-84964477feaa-c000.snappy.parquet</code>	parquet	February 17, 2021, 06:14:25 (UTC+01:00)	1.6 KB	Standard

Process log data (`log_data` directory).

We will perform ETL on the files in `log_data` directory to create the remaining two dimensional tables: time and users, as well as the songplays fact table.

This is what a single log file looks like:

```
{"artist":null,"auth":"Logged In","firstName":"Walter","gender":"M","itemInSession":  
{"artist":null,"auth":"Logged In","firstName":"Kaylee","gender":"F","itemInSession":  
{"artist":"Des'ree","auth":"Logged In","firstName":"Kaylee","gender":"F","itemInSess
```

We will proceed to read the data from the source S3 bucket, in the image below we can check the result of the task.

```
df = spark.read.json(log_data)
```

```
--- Reading log data from json files. --- 2021-02-17  
- Shows the scheme and some log data records...  
root  
|-- artist: string (nullable = true)  
|-- auth: string (nullable = true)  
|-- firstName: string (nullable = true)  
|-- gender: string (nullable = true)  
|-- itemInSession: long (nullable = true)  
|-- lastName: string (nullable = true)  
|-- length: double (nullable = true)  
|-- level: string (nullable = true)  
|-- location: string (nullable = true)  
|-- method: string (nullable = true)  
|-- page: string (nullable = true)  
|-- registration: double (nullable = true)  
|-- sessionId: long (nullable = true)  
|-- song: string (nullable = true)  
|-- status: long (nullable = true)  
|-- ts: long (nullable = true)  
|-- userAgent: string (nullable = true)  
|-- userId: string (nullable = true)  
-----  
|  artist|  auth|firstName|gender|itemInSession|lastName|  length|level|  location|method|  page|  registration|sessionId|  song|status|  ts|  
-----  
|  Harmonia|Logged In| Ryan|  M|  0| Smith|655.77751| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12|  583| Sehr kosmisch| 200|1542241826796|Mozilla/  
5.0 (X11...| 26|  
|The Prodigy|Logged In| Ryan|  M|  1| Smith|260.07465| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12|  583|The Big Gundown| 200|1542242481796|Mozilla/  
5.0 (X11...| 26|  
|  Train|Logged In| Ryan|  M|  2| Smith|205.45261| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12|  583|  Marry Me| 200|1542242741796|Mozilla/  
5.0 (X11...| 26|  
|  null|Logged In| Wyatt|  M|  0| Scott|  null| free|Eureka-Arcata-For...| GET|  Home|1.540872073796E12|  563|  null| 200|1542247071796|Mozilla/5  
.0 (Wind...| 9|  
|  null|Logged In| Austin|  M|  0| Rosales|  null| free|New York-Newark-J...| GET|  Home|1.541059521796E12|  521|  null| 200|1542252577796|Mozilla/5  
.0 (Wind...| 12|  
-----  
only showing top 5 rows
```

First we filter by actions for song plays.

```
df = df.where("page = 'NextSong'")
```

For the users table, we will extract and transform the log data using only the columns corresponding to the users table suggested in the project specifications, filtering the records that have value in the user id, sorting the records by user ID and without duplicates.

```
users_table = df.selectExpr("userId as user_id", "firstName as first_name", "las  
  .filter("user_id <> ''") \  
  .orderBy("user_id") \  
  .drop_duplicates()
```

```

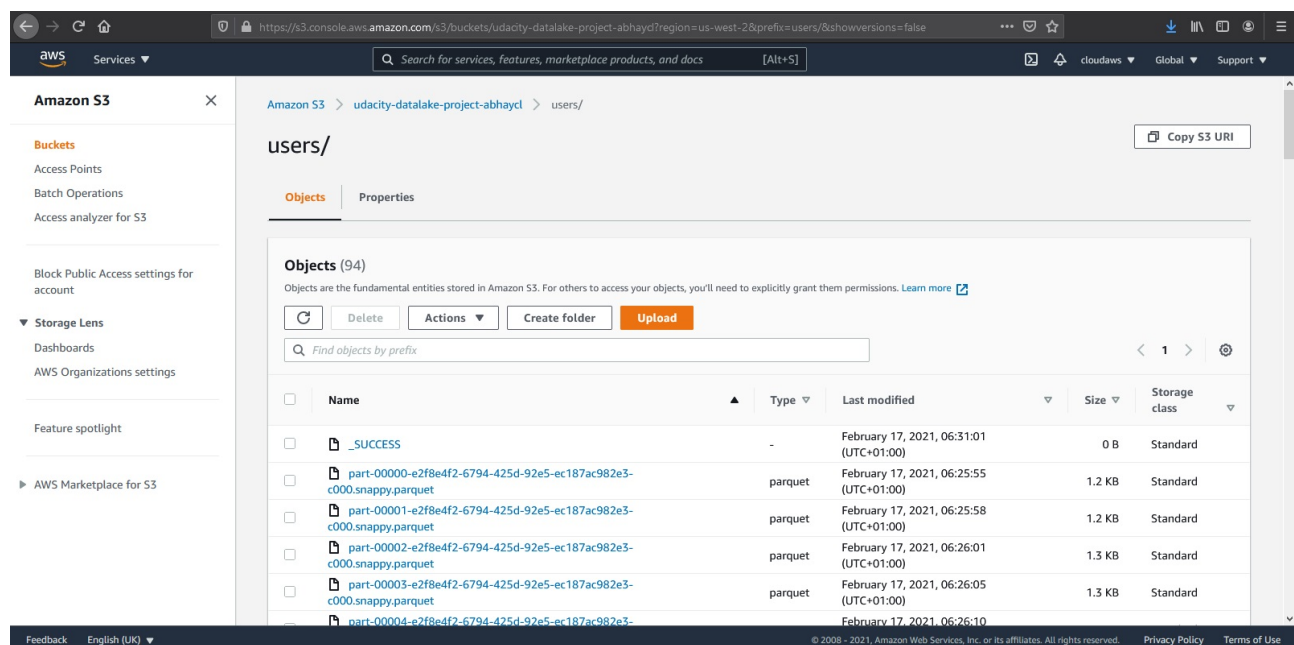
--- Extraction of columns to create users table. ---
- Shows the scheme and some users table records...
root
|-- user_id: string (nullable = true)
|-- first_name: string (nullable = true)
|-- last_name: string (nullable = true)
|-- gender: string (nullable = true)
|-- level: string (nullable = true)

+-----+-----+-----+-----+-----+
|user_id|first_name|last_name|gender|level|
+-----+-----+-----+-----+-----+
|    10|    Sylvie|    Cruz|    F| free|
|   100|    Adler|  Barrera|    M| free|
|   101|   Jayden|    Fox|    M| free|
|    11| Christian|  Porter|    F| free|
|    12|   Austin|  Rosales|    M| free|
+-----+-----+-----+-----+-----+
only showing top 5 rows

```

We proceed to load our data into our destination S3 bucket to write the users table in parquet files.

```
users_table.write.parquet(os.path.join(output_data, "users"), mode="overwrite")
```



Creation of the timestamp and datetime column from original timestamp column.

```

get_timestamp = udf(lambda x: datetime.fromtimestamp(int(x)/1000).strftime("%Y-%m-%d"))
df = df.withColumn("timestamp", get_timestamp("ts"))

```

```

get_datetime = udf(lambda x: datetime.fromtimestamp(int(x)/1000).strftime("%Y-%m-%d %H:%M:%S"))
df = df.withColumn("datetime", get_datetime("ts"))

```

```

--- Creation of timestamp and datetime columns. ---
+-----+
| artist| auth|firstName|gender|listenSession|lastName| length|level| location|method| page| registration|sessionId| song|status| ts|
+-----+
| Harmonia|Logged In| Ryan| M| 0| Smith|655.77751| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12| 583| Sehr Kosmisch| 200|1542241826796|"Moz
illa/5.0 (Xil...| 26|2018-11-15 00:30:26|2018-11-15| 1| Smith|260.07465| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12| 583| The Big Gundown| 200|1542242481796|"Moz
The Prodigy|Logged In| Ryan| M| 1| Smith|260.07465| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12| 583| The Big Gundown| 200|1542242481796|"Moz
illa/5.0 (Xil...| 26|2018-11-15 00:41:21|2018-11-15| 1| Smith|205.45261| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12| 583| Marry Me| 200|1542242741796|"Moz
Train|Logged In| Ryan| M| 2| Smith|205.45261| free|San Jose-Sunnyval...| PUT|NextSong|1.541016707796E12| 583| Marry Me| 200|1542242741796|"Moz
illa/5.0 (Xil...| 26|2018-11-15 00:45:41|2018-11-15| 1| Gonzalez|218.06975| free|Houston-The Woodl...| PUT|NextSong|1.540492941796E12| 597| Blackbird| 200|1542253449796|"Moz
Sony Wonder|Logged In| Samuel| M| 0| Gonzalez|218.06975| free|Houston-The Woodl...| PUT|NextSong|1.540492941796E12| 597| Blackbird| 200|1542253449796|"Moz
illa/5.0 (Mac...| 61|2018-11-15 03:44:09|2018-11-15| 1| Levine|289.38404| paid|Portland-South Po...| PUT|NextSong|1.540794356796E12| 602|Best Of Both Worl...| 200|1542260935796|"Moz
Van Halen|Logged In| Tegan| F| 2| Levine|289.38404| paid|Portland-South Po...| PUT|NextSong|1.540794356796E12| 602|Best Of Both Worl...| 200|1542260935796|"Moz
illa/5.0 (Mac...| 80|2018-11-15 05:48:55|2018-11-15| 1|
+-----+
only showing top 5 rows

```

For the time table, we will extract and transform the log data using only the columns corresponding to the time table suggested in the project specifications, sorting the records by start time and without duplicates.

```

time_table = df.select(col("timestamp").alias("start_time"), hour("timestamp").a
                .orderBy("start_time") \
                .drop_duplicates()

```

```

--- Extraction of columns to create time table. ---
- Shows the scheme and some time table records...
root
 |-- start_time: string (nullable = true)
 |-- hour: integer (nullable = true)
 |-- day: integer (nullable = true)
 |-- week: integer (nullable = true)
 |-- month: integer (nullable = true)
 |-- year: integer (nullable = true)
 |-- weekday: string (nullable = true)

+-----+-----+-----+-----+-----+-----+-----+
| start_time|hour|day|week|month|year|weekday|
+-----+-----+-----+-----+-----+-----+-----+
|2018-11-01 21:01:46| 21| 1| 44| 11|2018| 4|
|2018-11-01 21:05:52| 21| 1| 44| 11|2018| 4|
|2018-11-01 21:08:16| 21| 1| 44| 11|2018| 4|
|2018-11-01 21:11:13| 21| 1| 44| 11|2018| 4|
|2018-11-01 21:17:33| 21| 1| 44| 11|2018| 4|
+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows

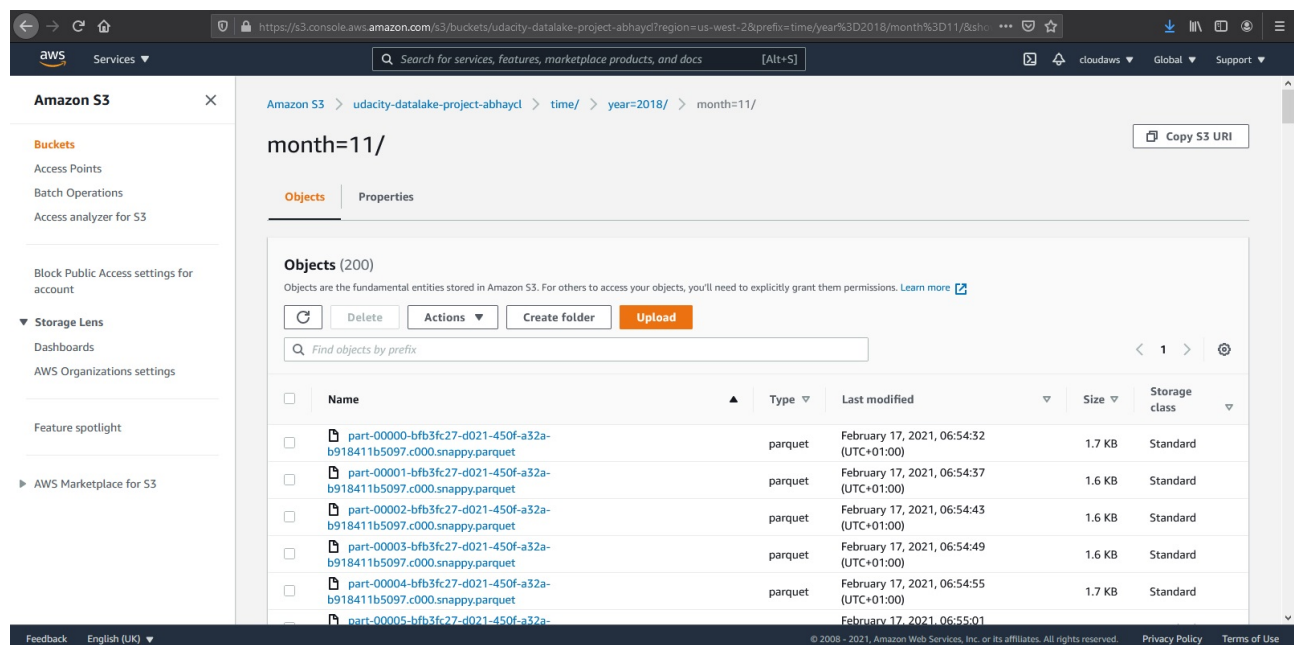
```

We proceed to load our data into our destination S3 bucket to write the song table in parquet files partitioned by year and month.

```

time_table.write.partitionBy("year", "month").parquet(os.path.join(output_da

```



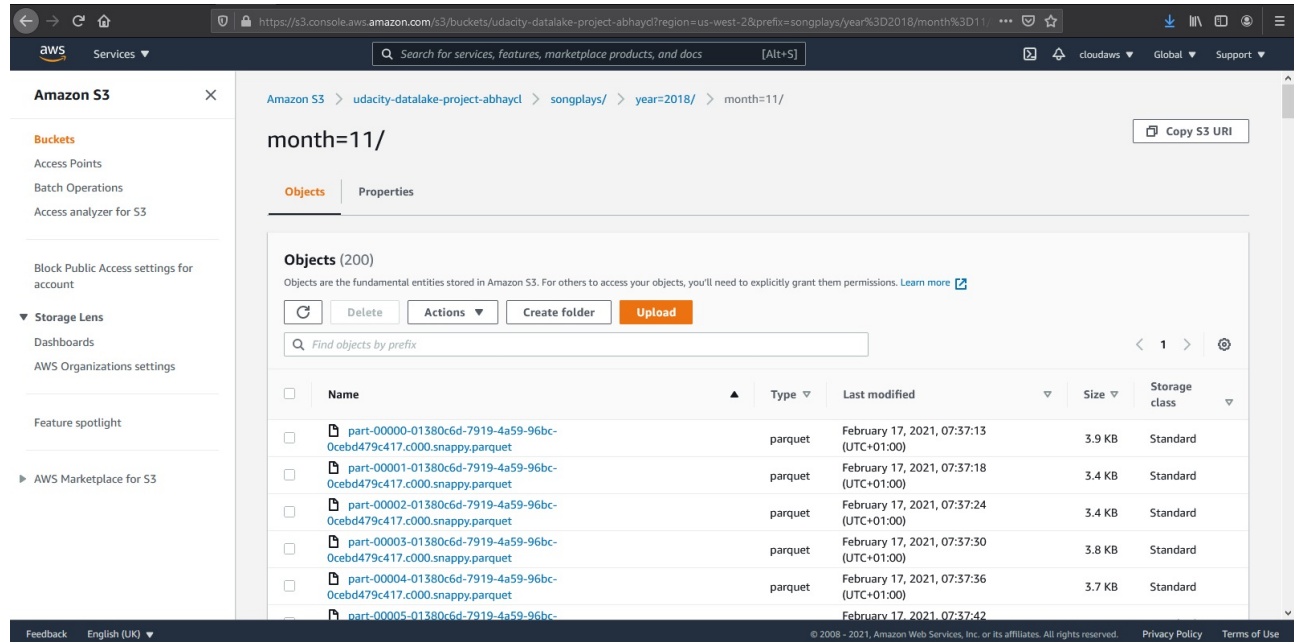
```
song_df = spark.read.json(os.path.join(input_data, "song_data/**/*.json")).se
```

```
songplays_table = df.join(song_df, \
    (df.song == song_df.title) & (df.artist == song_df.artist_name) & (df.length
    .select(df.timestamp.alias("start_time"), df.userId.alias("user_id"), df.lev
    .orderBy("start_time", "user_id") \
    .withColumn("songplay_id", F.monotonically_increasing_id()))
```

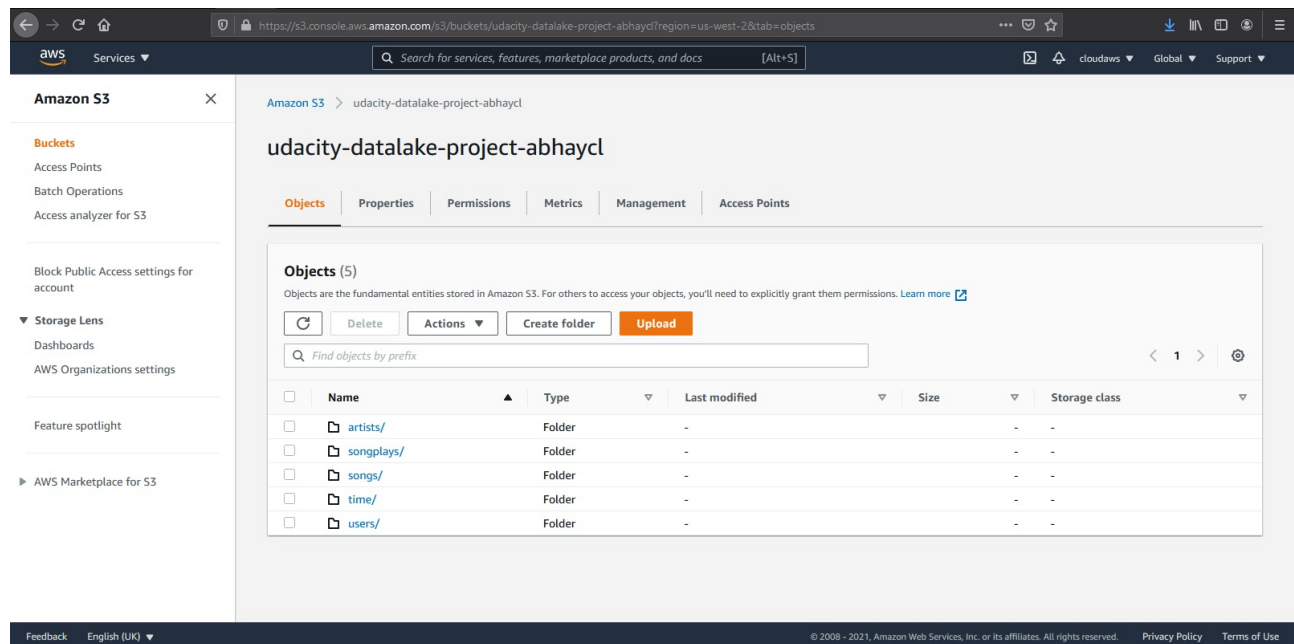
[illegible]

We proceed to load our data into our destination S3 bucket to write the songplays table in parquet files partitioned by year and month.

```
songplays_table.write.partitionBy("year", "month").parquet(os.path.join(outp
```



In our target bucket S3 we will be able to check all saved data corresponding to the five defined tables.



Conclusion.

As the data engineer of the music streaming startup. We read our data from a source S3 bucket, with our spark cluster we extracted and transformed the data into five tables according to the project specifications and deposited them in our destination S3 bucket for further analysis or data processing.