# Project

# Music Data Analysis

## Data Ingestion and Initial Validation

- ### Generate_mob_data

**Python script is used to generate random data**
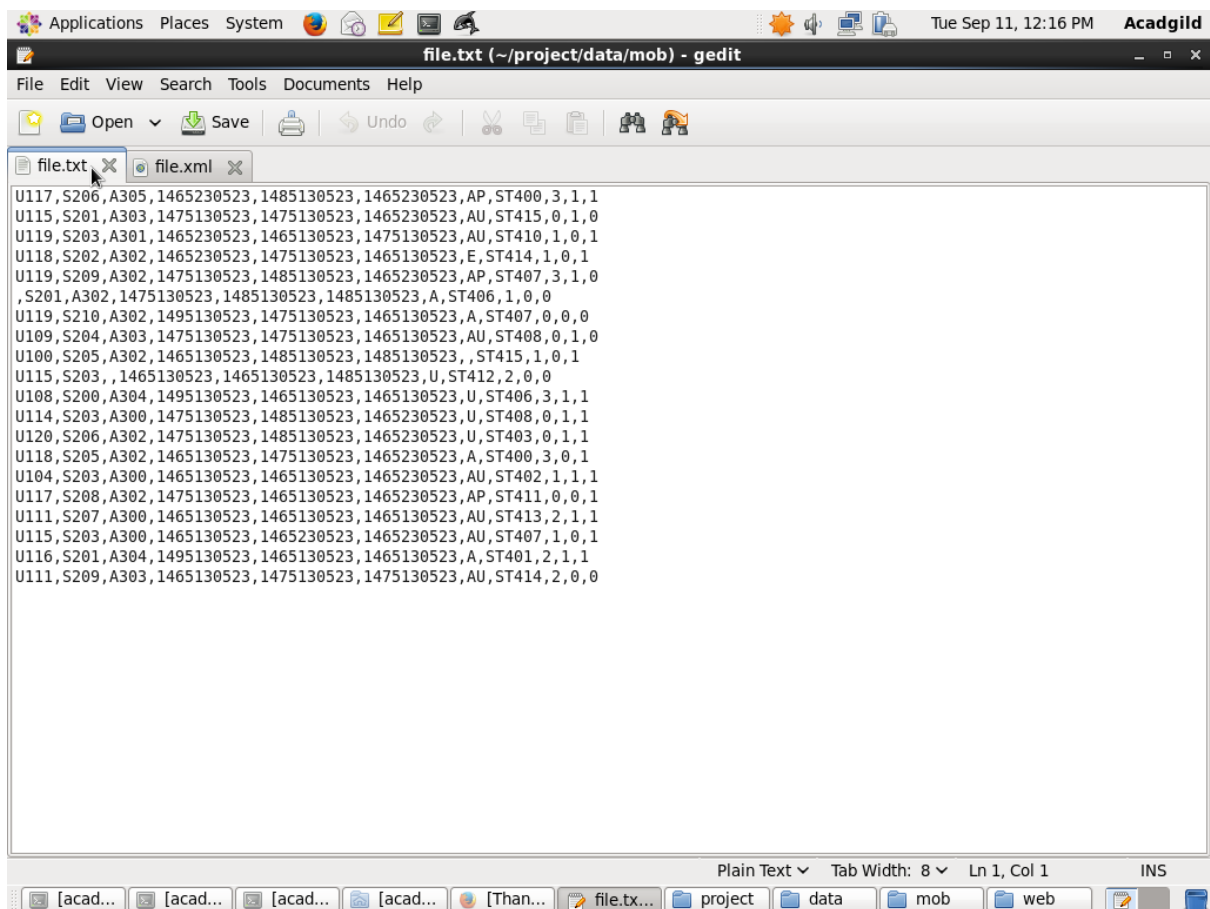
**Command:**

**python /home/acadgild/project/scripts/generate_mob_data.py**

**Code: (with comments)**

generate_mob_data.p
y

**Screen:**

- **Generate_web_data**

**Python script is used to generate random data**

**Command:**

python /home/acadgild/project/scripts/generate_web_data.py

**Code:**



generate_web_data.py

**Screen:**

- **Starting Hadoop**

sh /home/acadgild/project/scripts/start-daemons.sh

start-daemons.sh

- **Populate-Lookup**

sh /home/acadgild/project/scripts/populate_lookup.sh

populate-lookup.sh



user-artist.hql

- **Data Formating:**

python /home/acadgild/project/scripts/data_formating.sh

dataformatting.sh



dataformatting.pig
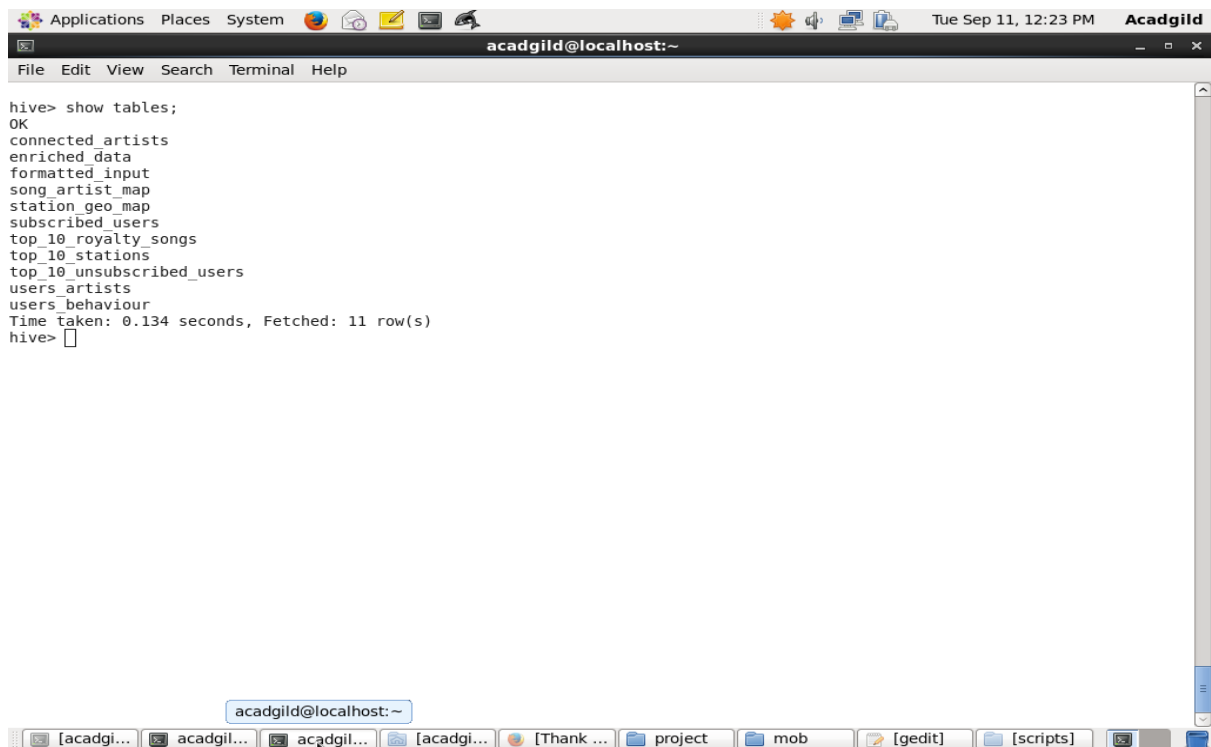


formatted_hive_load.
hql

Terminal screenshot content:

```
subscribed_users
top_10_royalty_songs
top_10_stations
top_10_unsubscribed_users
users_artists
users_behaviour
Time taken: 0.134 seconds, Fetched: 11 row(s)
hive> select * from enriched_data
    > LIMIT 10;
OK
U104    S202    A302    1462863262    1465490556    1465490556    A      ST410   1   1   1   1   fail
U117    S203    A303    1465130523    1465130523    1485130523    A      ST400   1   1   1   1   fail
U109    S204    A304    1495130523    1485130523    1465130523    NULL   ST415   2   0   1   1   fail
U115    S204    A304    1494297562    1494297562    1468094889    J      ST413   0   1   1   1   fail
U100    S206    A302    1468094889    1465490556    1465490556    E      ST409   1   1   1   1   fail
U107    S206    A302    1462863262    1494297562    1462863262    NULL   ST415   2   0   1   1   fail
U102    S206    A302    1494297562    1462863262    1494297562    A      ST400   1   1   1   1   fail
U104    S207    A303    1475130523    1475130523    1465230523    A      ST411   1   1   1   1   fail
U109    S207    A303    1465230523    1475130523    1465230523    A      ST405   0   1   1   1   fail
U102    S208    A304    1475130523    1465130523    1475130523    E      ST414   2   1   1   1   fail
Time taken: 0.715 seconds, Fetched: 10 row(s)
hive> select * from formated_input
    > LIMIT 10;
FAILED: SemanticException [Error 10001]: Line 1:14 Table not found 'formated_input'
hive> select * from formatted_input
    > LIMIT 10;
OK
U106    S210    A305    1475130523    1475130523    1465130523    E     ST413   0   1   1   1
U104    S207    A304    1475130523    1475130523    1465230523    A     ST411   1   1   1   1
U102    S208    A302    1475130523    1465130523    1475130523    AU    ST414   2   1   1   1
U104    S206    A304    1475130523    1465130523    1475130523    E     ST410   0   0   1   1
U118    S204    A302    1465130523    1465130523    1485130523    U     ST406   2   1   0   1
        S209    A300    1475130523    1465130523    1465130523    A     ST409   2   1   0   1
U109    S204    A304    1495130523    1485130523    1465130523    U     ST415   2   0   1   1
U107    S209    A303    1495130523    1465130523    1475130523    E     ST401   1   0   1   1
U116    S208    A302    1495130523    1465130523    1475130523          ST407   1   0   1   1
U101    S200            1475130523    1465130523    1475130523    E     ST402   2   1   0   1
Time taken: 0.232 seconds, Fetched: 10 row(s)
hive>
```

- **Hive lookup:**

**Command:**

**Hive -f /home/acadgild/project/scripts/create_hive_hbase_lookup.hql**

**Code:**



create_hive_hbase_lo
okup.hql

**Screen:**

acadgild@localhost:~

File  Edit  View  Search  Terminal  Help

```
hive> show tables;
OK
connected_artists
enriched_data
formatted_input
song_artist_map
station_geo_map
subscribed_users
top_10_royalty_songs
top_10_stations
top_10_unsubscribed_users
users_artists
users_behaviour
Time taken: 0.134 seconds, Fetched: 11 row(s)
hive>
```

acadgild@localhost:~

[acadgi...]  acadgil...  acadgil...  [acadgi...  [Thank ...  project  mob  [gedit]  [scripts]

## Data Enrichment

## **Command:**

**sh /home/acadgild/project/scripts/data_enrichment.sh**

## **Code:**

data_enrichment.sh      data_enrichment.hql

## **Screen:**

```
                              acadgild@localhost:~                            _  □  ✕

File  Edit  View  Search  Terminal  Help

hive> show tables;
OK
connected_artists
enriched_data
formatted_input
song_artist_map
station_geo_map
subscribed_users
top_10_royalty_songs
top_10_stations
top_10_unsubscribed_users
users_artists
users_behaviour
Time taken: 0.134 seconds, Fetched: 11 row(s)
hive> select * from enriched_data
    > LIMIT 10;
OK
U104    S202    A302    1462863262    1465490556    1465490556    A     ST410  1  1  1  1  fail
U117    S203    A303    1465130523    1465130523    1485130523    A     ST400  1  1  1  1  fail
U109    S204    A304    1495130523    1485130523    1465130523    NULL  ST415  2  0  1  1  fail
U115    S204    A304    1494297562    1494297562    1468094889    J     ST413  0  1  1  1  fail
U100    S206    A302    1468094889    1465490556    1465490556    E     ST409  1  1  1  1  fail
U107    S206    A302    1462863262    1494297562    1462863262    NULL  ST415  2  0  1  1  fail
U102    S206    A302    1494297562    1462863262    1494297562    A     ST400  1  1  1  1  fail
U104    S207    A303    1475130523    1475130523    1465230523    A     ST411  1  1  1  1  fail
U109    S207    A303    1465230523    1475130523    1465230523    A     ST405  0  1  1  1  fail
U102    S208    A304    1475130523    1465130523    1475130523    E     ST414  2  1  1  1  fail
Time taken: 0.715 seconds, Fetched: 10 row(s)
hive> []
```

[acadgi...] [ acadgil...] [ acadgil...] [ [acadgi...] [Thank ...] [ project] [ mob] [ [gedit]] [ [scripts]] [▣] [■]

## Data Analysis

1. Determine top 10 station_id(s) where maximum number of songs were played, which were liked by unique users.
2. Determine total duration of songs played by each type of user, where type of user can be **'subscribed'** or **'unsubscribed'**. An unsubscribed user is the one whose record is either not present in **Subscribed_users** lookup table or has *subscription_end_date* earlier than the *timestamp* of the song played by him.
3. Determine top 10 connected artists. Connected artists are those whose songs are most listened by the unique users who follow them.
4. Determine top 10 songs who have generated the maximum revenue. Royalty applies to a song only if it was *liked* or was *completed successfully* or both.
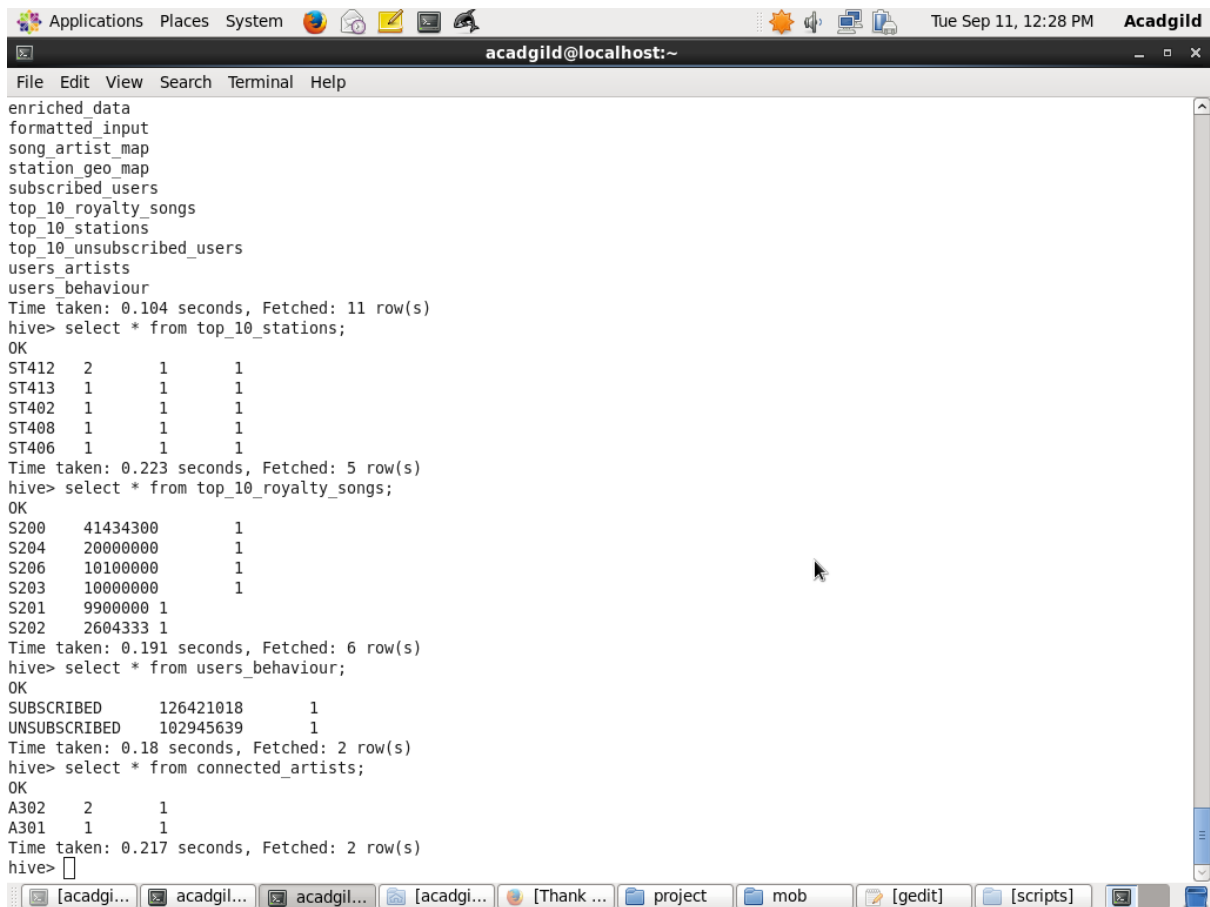
## Command:

**sh /home/acadgild/project/scripts/data_analysis.sh**

## Code:

data_analysis.sh        data_analysis.hql        data_export.sh        create_schema.sql

## Screen:

## Hive:



## Mysql create schema:

```
mysql>
mysql> show tables;
+---------------------------+
| Tables_in_project         |
+---------------------------+
| connected_artists         |
| top_10_royalty_songs      |
| top_10_stations           |
| top_10_unsubscribed_users |
| users_behaviour           |
+---------------------------+
5 rows in set (0.07 sec)
```

## Mysql Analysis:

File  Edit  View  Search  Terminal  Help

```
mysql> select * from connected_artists;
+-----------+-------------+
| artist_id | user_count  |
+-----------+-------------+
| A302      |           2 |
| A301      |           1 |
+-----------+-------------+
2 rows in set (0.01 sec)

mysql> select * from top_10_stations;
+------------+----------------------------+--------------------+
| station_id | total_distinct_songs_played | distinct_user_count |
+------------+----------------------------+--------------------+
| ST412      |                          2 |                  1 |
| ST413      |                          1 |                  1 |
| ST402      |                          1 |                  1 |
| ST408      |                          1 |                  1 |
| ST406      |                          1 |                  1 |
+------------+----------------------------+--------------------+
5 rows in set (0.00 sec)

mysql> select * from users_behavious;
ERROR 1146 (42S02): Table 'project.users_behavious' doesn't exist
mysql> select * from users_behaviour;
+--------------+-----------+
| user_type    | duration  |
+--------------+-----------+
| SUBSCRIBED   | 126421018 |
| UNSUBSCRIBED | 102945639 |
+--------------+-----------+
2 rows in set (0.02 sec)

mysql> select * from top_10_royalty_songs;
+---------+----------+
| song_id | duration |
+---------+----------+
| S200    | 41434300 |
| S204    | 20000000 |
| S206    | 10100000 |
```

[acadgi...]  acadgil...  [acadgi...]  [acadgi...]  [Thank ...]  project  mob  [gedit]  [scripts]

```
mysql> select * from top_10_royalty_songs;
+---------+----------+
| song_id | duration |
+---------+----------+
| S200    | 41434300 |
| S204    | 20000000 |
| S206    | 10100000 |
| S203    | 10000000 |
| S201    |  9900000 |
| S202    |  2604333 |
+---------+----------+
6 rows in set (0.00 sec)

mysql>
```

## Job scheduling:

## Command:

- **Open cron tab -e**
- **Add the code for scheduling:**
    **\* \*/3 \* \* \* /home/acadgild/project/scripts/wrapper.sh**

## Code:

wrapper.sh

<span style="color:red">**Screen:**</span>

```
File  Edit  View  Search  Terminal  Help
[acadgild@localhost ~]$ sudo crontab -e
[sudo] password for acadgild:
no crontab for root - using an empty one
crontab: installing new crontab
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost ~]$
```