Assignment Submission-Final Project (Music Data Analysis)

1. Project Description:

A leading music-catering company is planning to analyze large amount of data received from varieties of sources, namely mobile app and website to track the behavior of users, classify users, calculate royalties associated with the song and make appropriate business strategies. The file server receives data files periodically after every 3 hours.

1.1 Fields present in the data files

Data files contain below fields.

Column Name/Field Name	Column Description/Field Description			
User_id	Unique identifier of every user			
Song_id	Unique identifier of every song			
Artist_id	Unique identifier of the lead artist of the song			
Timestamp	Timestamp when the record was generated			
Start_ts	Start timestamp when the song started to play			
End_ts	End timestamp when the song was stopped			
Geo_cd	Can be 'A' for USA region, 'AP' for asia pacific region,'J' for Japan region, 'E' for europe and 'AU' for australia region			
Station_id	Unique identifier of the station from where the song was played			
Song_end_type	How the song was terminated. 0 means completed successfully 1 means song was skipped 2 means song was paused 3 means other type of failure like device issue, network error etc.			
Like	0 means song was not likedsong was played 1 means song was liked			
Dislike	0 means song was not disliked 1 means song was disliked			

1.2 LookUp Tables

There is some existing look up tables present in **NoSQL** databases. They play an important role in data enrichment and analysis.

Table Name	Description		
Station_Geo_Map	Contains mapping of a geo_cd with station_id		
Subscribed_Users	Contains user_id, subscription_start_date and subscription_end_date. Contains details only for subscribed users		
Song_Artist_Map	Contains mapping of song_id with artist_id alongwith royalty associated with each play of the song		
User_Artist_Map	Contains an array of artist_id(s) followed by a user_id		

1.3 Data Ingestion and Initial Validation

Below is the link for datasets.

https://drive.google.com/drive/folders/OB P3pWagdIrrMjJGVINsS

UEtbG8

- 1. Data coming from web applications reside in /data/web and has xml format.
- 2. Data coming from mobile applications reside in /data/mob and has csv format.
- 3. Data files come every 3 hours.
- 4. All the timestamp fields in data coming from web application is of the format YYYY-MM-DD HH:MM:SS.
- 5. All the timestamp fields in data coming from mobile application is a long integer interpreted as UNIX timestamps.
- 6. Finally, all timestamps must have the format of a long integer to be interpreted as UNIX timestamps.
- 7. If both like and dislike are 1, consider that record to be invalid.
- 8. If any of the fields from User_id, Song_id, Timestamp, Start_ts, End_ts, Geo_cd is NULL or absent, consider that record to be invalid.
- 9. If Song_end_type is NULL or absent, treat it to be 3.
- 10. Create a temporary identifier for all the data files received in the last 3 hours (may be an integer batch_id which is auto incremented or a string obtained after combining current date and current hour, to keep track of valid and invalid records per batch).

1.4 Data Enrichment

- 1. If any of like or dislike is NULL or absent, consider it as 0.
- 2. If fields like Geo_cd and Artist_id are NULL or absent, consult the lookup tables for fields Station_id and Song_id respectively to get the values of Geo_cd and Artist_id.
- 3. If corresponding lookup entry is not found, consider that record to be invalid.

NULL or absent field	Look up field	Look up table (Table from which record can be updated)
Geo_cd	Station_id	Station_Geo_Map
Artist_id	Song_id	Song_Artist_Map

1.5 Data Analysis

It is not only the data which is important, rather it is the insight it can be used to generate important. Once we have made the data ready for analysis, we have to perform below analysis on a daily basis.

- 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 presents 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 that have generated the maximum revenue. Royalty applies to a song only if it was liked or was completed successfully or both.
- 5. Determine top 10 unsubscribed users who listened to the songs for the longest duration.

PROJECT IMPLEMENTATION

1. Data creation:

We have generated data through python scripts. Those python scripts are:

```
generate_web_data.py
generate_mob_data.py
```

Data coming from web applications reside in /home/acadgild/examples/music/data/web and has xml format and that coming from mobile applications reside in /home/acadgild/examples/music/data/mob and has text format.

The batch file "music_project_master.sh" does data creation through python scripts. Please find below script which is part of music_project_master.sh:

So here it will first remove web and mob directories, if they are present already inside directory

/home/acadgild/examples/music/data.

Then it will create web and mob directories inside directory

/home/acadgild/examples/music/data.

```
[acadgild@localhost music]$ ./music_project_master.sh
Preparing to execute python scripts to generate data...
rm: cannot remove `/home/acadgild/examples/music/data/web': No such file or directory
rm: cannot remove `/home/acadgild/examples/music/data/mob': No such file or directory
Data Generated Successfully !

Ctarting the degrees
```

2. Start Hadoop Daemons:

Created a batch file "start-daemon.sh". Through this batch file, these daemons are started. Please find below:

```
***
#!/bin/bash
rm -r /home/acadgild/examples/music/logs
mkdir -p /home/acadgild/examples/music/logs
if [ -f "/home/acadgild/examples/music/logs/current-batch.txt" ]
 echo "Batch File Found!"
else
 echo -n "1" > "/home/acadgild/examples/music/logs/current-
batch.txt"
fi
chmod 775 /home/acadgild/examples/music/logs/current-batch.txt
echo "After chmod"
batchid=`cat /home/acadgild/examples/music/logs/current-batch.txt`
echo "After batchid-->> "$batchid
LOGFILE=/home/acadgild/examples/music/logs/log batch $batchid
echo "Starting daemons" >> $LOGFILE
start-all.sh
start-hbase.sh
mr-jobhistory-daemon.sh start historyserver
cat /home/acadgild/examples/music/logs/current-batch.txt
```

Here, it will first remove logs directory, if they are present inside directory /home/acadgild/examples/music/.

Then it will create logs directory inside directory /home/acadgild/examples/music/.

After this, it will search for **current-batch.txt** file inside directory: **/home/acadgild/examples/music/logs.**

If it is present, then message will be present as "Batch File Found" else it will create current-batch.txt file inside directory: /home/acadgild/examples/music/logs with content as '1'.

After this required permissions would be given for this file.

"batched" would be content of **current-batch.txt** file. i.e **1**. After this, log_batch_1 file as Logfile would be created inside directory **/home/acadgild/examples/music/logs/.**

"start-daemon.sh" batch file will start though music project master.sh batch file.

3. Populate Look up tables in HBase:

By using the "populate-lookup.sh" script, we will create below lookup tables in HBase. These tables we are using for Data formatting, Data enrichment and Analysis stage.

Table Name	Description		
Station_Geo_Map	Contains mapping of a geo_cd with station_id		
Subscribed_Users	Contains user_id, subscription_start_date and subscription_end_date. Contains details only for subscribed users		
Song_Artist_Map	Contains mapping of song_id with artist_id alongwith royalty associated with each play of the song		
User_Artist_Map	Contains an array of artist_id(s) followed by a user_id		

The "populate-lookup.sh" shell script creates the above lookup tables in the HBase and populate the data into the lookup tables from the dataset files.

```
*****
#!/bin/bash
batchid=`cat /home/acadgild/examples/music/logs/current-batch.txt`
LOGFILE=/home/acadgild/examples/music/logs/log batch $batchid
echo "Creating LookUp Tables" >> $LOGFILE
echo "disable 'station-geo-map'" | hbase shell echo "drop 'station-
geo-map'" | hbase shell
echo "disable 'subscribed-users'" | hbase shell
echo "drop 'subscribed-users'" | hbase shell
echo "disable 'song-artist-map'" | hbase shell
echo "drop 'song-artist-map'" | hbase shell
echo "create 'station-geo-map', 'geo'" | hbase shell
echo "create 'subscribed-users', 'subscn'" | hbase shell
echo "create 'song-artist-map', 'artist'" | hbase shell
echo "Populating LookUp Tables" >> $LOGFILE
file="/home/acadgild/examples/music/lookupfiles/stn-geocd.txt"
while IFS= read -r line
do
 stnid=`echo $line | cut -d',' -f1`
 geocd=`echo $line | cut -d',' -f2`
 echo "put 'station-geo-map', '$stnid', 'geo:geo cd', '$geocd'" |
hbase shell
done <"$file"</pre>
file="/home/acadgild/examples/music/lookupfiles/song-artist.txt"
while IFS= read -r line
do
 songid=`echo $line | cut -d',' -f1`
 artistid=`echo $line | cut -d',' -f2`
 echo "put 'song-artist-map', '$songid', 'artist:artistid',
'$artistid'" | hbase shell
done <"$file"</pre>
```

file="/home/acadgild/examples/music/lookupfiles/user-subscn.txt"
while IFS= read -r line

Below screen shots shows the tables creation and population of the data in the HBase. Here we are executing **populate-lookup.sh** via music_project_master.sh batch file.

```
[acadgild@localhost music]$ ./music_project_master.sh
Preparing to execute python scripts to generate data...
Data Generated Successfully !
 Starting the daemons....
12514 Jps
5095 DataNode
5257 SecondaryNameNode
6377 JobHistoryServer
  5257 SecondaryNameNode
6377 JobHistoryServer
5001 NameNode
  5484 ResourceManager
5583 NodeManager
11985 Main
11985 Main
6131 HQuorumPer
7380 Main
6196 HMaster
7576 RunJar
6297 HRegionServer
All hadoop daemons started !
Upload the look up tables now in Hbase...
2018-11-25 22:01:21,718 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
SLF4J: class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
 .class]
SLF4J: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/i
mpl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type lorg.slf4j.impl.Log4jLoggerFactory]
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
  disable 'station-geo-map'
0 row(s) in 6.0110 seconds
 2018-11-25 22:02:19,615 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder
 SLF41: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF41: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF41: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF41: Actual binding is of type [org.slf4j.impl.log4jloggerFactory]
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
  drop 'station-geo-map'
0 row(s) in 5.0480 seconds
2018-11-25 22:03:16,646 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit-RETURN>" to leave the HBase Shell
Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
  disable 'subscribed-users'
0 row(s) in 5.7790 seconds
  2018-11-25 22:04:14,135 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
```

```
subscribed-users'
    row(s) in 3.9400 seconds
 2018-11-25 22:05:08,393 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder
.class]
SLF41: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/i
mpl/StaticLoggerBinder.class]
SLF41: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF43: Actual binding is of type [org.slf4j.impl.log4jLoggerFactory]
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
disable 'song-artist-map'
0 row(s) in 5.0760 seconds
 2018-11-25 22:06:01,334 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder
SLF4J: Found binding in []ar:file:/home/acadgild/instatt/hodse/.bdop/.class]
SLF4J: Found binding in [jar:file:/home/acadgild/instatl/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
HBase Shell; enter 'help=RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
drop 'song-artist-map'
0 row(s) in 4.2220 seconds
2018-11-25 22:06:55,326 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
classes where applicable
 create 'station-geo-map', 'geo'
 0 row(s) in 4.9870 seconds
 Hbase::Table - station-geo-map
2018-11-25 22:07:51,435 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder
 .class]
SLF4J: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
mpt/staticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.log4jLoggerFactory]
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
 create 'subscribed-users', 'subscn'
    row(s) in 4.1690 seconds
  Hbase::Table - subscribed-users
 2018-11-25 22:08:46,537 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
 classes where applicable
SLF4J: Class path contains multiple SLF4J bindings.
 SLF4J: Found binding in [jar:file:/home/acadgild/install/hbase/hbase-1.2.6/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder
 .class]
.class|
SLF4J: Found binding in [jar:file:/home/acadgild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class|
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
HBase Shell: enter 'helpeRETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.6. glukhown Mon May 29.03:35:32 CDT 2017
 Version 1.2.6, rUnknown, Mon May 29 02:25:32 CDT 2017
```

create 'song-artist-map', 'artist' 0 row(s) in 3.8860 seconds

```
hbase(main):002:0> scan 'song-artist-map
                                            COLUMN+CELL
 5200
                                            column=artist:artistid, timestamp=1532277052808, value=A300
 5201
                                            column=artist:artistid, timestamp=1532277063975, value=A301
                                           column=artist:artistid, timestamp=1532277074927, value=A302
column=artist:artistid, timestamp=1532277085940, value=A303
 5202
 5203
                                           column=artist:artistid, timestamp=1532277096508, value=A304 column=artist:artistid, timestamp=1532277107380, value=A301
 S204
 5205
                                            column=artist:artistid, timestamp=1532277117916, value=A302
 S206
 S207
                                            column=artist:artistid, timestamp=1532277128708, value=A303
 5208
                                            column=artist:artistid, timestamp=1532277139626, value=A304
9 row(s) in 0.2440 seconds
```

```
hbase(main):003:0> scan 'station-geo-map
                                                      COLUMN+CELL
ROW
                                                      column=geo:geo_cd, timestamp=1532276901571, value=Acolumn=geo:geo_cd, timestamp=1532276912108, value=AU
 ST400
 ST401
 ST402
                                                       column=geo:geo_cd, timestamp=1532276922831, value=AP
                                                      column=geo:geo_cd, timestamp=1532276933380, value=J
column=geo:geo_cd, timestamp=1532276944269, value=E
 ST403
 ST404
                                                       column=geo:geo_cd, timestamp=1532276954714, value=A
 ST405
                                                      column=geo:geo_cd, timestamp=1532276966054, value=AUcolumn=geo:geo_cd, timestamp=1532276976538, value=AP
 ST406
 ST407
                                                       column=geo:geo_cd, timestamp=1532276987193, value=E
 ST408
                                                      column=geo:geo_cd, timestamp=1532276998216, value=E
column=geo:geo_cd, timestamp=1532277009161, value=A
column=geo:geo_cd, timestamp=1532277020083, value=A
 ST409
 ST410
 ST411
                                                      column=geo:geo_cd, timestamp=1532277030853, value=AP
column=geo:geo_cd, timestamp=1532277041902, value=J
 ST412
 ST413
14 row(s) in 0.1260 seconds
```

```
hbase(main):004:0> scan 'subscribed-users
                                                                                                                                                                                                                                                                                          COLUMN+CELL
  ROW
                                                                                                                                                                                                                                                                                       column=subscn:enddt, timestamp=1532277161513, value=1465130523
column=subscn:startdt, timestamp=1532277150572, value=1465230523
column=subscn:enddt, timestamp=1532277183558, value=1475130523
column=subscn:startdt, timestamp=1532277172591, value=1465230523
column=subscn:enddt, timestamp=1532277194398, value=1475130523
column=subscn:enddt, timestamp=1532277194398, value=1465230523
      U100
      U100
      U101
      U101
      U102
                                                                                                                                                                                                                                                                                   column=subscn:enddt, timestamp=1532277194398, value=1475130523 column=subscn:startdt, timestamp=1532277194398, value=1475130523 column=subscn:enddt, timestamp=1532277226044, value=1475130523 column=subscn:startdt, timestamp=1532277215490, value=1465230523 column=subscn:enddt, timestamp=1532277248517, value=1475130523 column=subscn:enddt, timestamp=1532277237099, value=1465230523 column=subscn:enddt, timestamp=1532277270534, value=1475130523 column=subscn:enddt, timestamp=1532277259547, value=1465230523 column=subscn:enddt, timestamp=1532277292198, value=1485130523 column=subscn:enddt, timestamp=1532277281420, value=1465230523 column=subscn:enddt, timestamp=1532277313425, value=1465230523 column=subscn:enddt, timestamp=1532277332798, value=1465230523 column=subscn:enddt, timestamp=1532277332798, value=1465230523 column=subscn:enddt, timestamp=15322773323818, value=1465230523 column=subscn:enddt, timestamp=1532277356273, value=1475130523 column=subscn:enddt, timestamp=1532277345393, value=1465230523 column=subscn:enddt, timestamp=1532277389179, value=1475130523 column=subscn:enddt, timestamp=1532277389179, value=1465230523 column=subscn:enddt, timestamp=1532277421027, value=1475130523 column=subscn:enddt, timestamp=1532277421027, value=145530523 column=subscn:enddt, timestamp=1532277421027, value=145530523 column=subscn:enddt, timestamp=1532277421027, value=1465230523 column=subscn:enddt, timestamp=1532277421027, value=1465230523 column=subscn:enddt, timestamp=1532277421027, value=1465230523 column=subscn:enddt, timest
        U102
        U103
      U103
      U104
        U104
        U105
        U105
      U106
        U106
        U107
        U107
        U108
        U108
        U109
        U109
        U110
        U110
        U111
        U112
        U112
        U113
     14 row(s) in 0.1340 seconds
```

With the help of Hbase storage handler & SerDe properties, we are creating the hive external tables by matching the columns of Hbase tables to hive tables.

data_enrichment_filtering_schema.sh script will run the "create_hive_hbase_lookup.hql" which will create the HIVE external tables.

```
data_enrichment_filtering_schema.sh script :
**********************
#!/bin/bash
batchid=`cat /home/acadgild/examples/music/logs/current-batch.txt`
LOGFILE=/home/acadgild/examples/music/logs/log batch $batchid
echo "Creating hive tables on top of hbase tables for data
enrichment and filtering..." >> $LOGFILE
hive -f /home/acadgild/examples/music/ create hive hbase lookup.hql
"create_hive_hbase_lookup.hql" script :
CREATE DATABASE IF NOT EXISTS project;
USE project;
create external table if not exists station geo map
station id String,
geo cd string
)
STORED BY 'org.apache.hadoop.hive.hbase.HBaseStorageHandler'
with serdeproperties
("hbase.columns.mapping"=":key,geo:geo cd")
tblproperties("hbase.table.name"="station-geo-map");
create external table if not exists subscribed users
(
user id STRING,
subscn_start_dt STRING,
subscn end dt STRING
)
STORED BY 'org.apache.hadoop.hive.hbase.HBaseStorageHandler'
with serdeproperties
("hbase.columns.mapping"=":key,subscn:startdt,subscn:enddt")
tblproperties("hbase.table.name"="subscribed-users");
create external table if not exists song artist map
song id STRING,
artist id STRING)
```

STORED BY 'org.apache.hadoop.hive.hbase.HBaseStorageHandler'

We are running **data_enrichment_filtering_schema.sh** script through the execution of **music_project_master.sh** script.

The below screenshot we can see tables are getting created in hive by running the

"data enrichement filtering schema.sh file".

```
hive> use project;
OK
Time taken: 5.868 seconds
hive> show tables;
OK
song_artist_map
station_geo_map
subscribed_users
users_artists
Time taken: 0.329 seconds, Fetched: 4 row(s)
```

```
hive> select * from song_artist_map;
OK
5200
        A300
S201
        A301
S202
        A302
S203
        A303
S204
        A304
S205
        A301
S206
        A302
S207
        A303
S208
        A304
Time taken: 4.118 seconds, Fetched: 9 row(s)
```

```
hive> select * from subscribed users;
OK
U100
        1465230523
                        1465130523
U101
        1465230523
                        1475130523
U102
        1465230523
                        1475130523
U103
        1465230523
                        1475130523
U104
        1465230523
                        1475130523
U105
        1465230523
                        1475130523
U106
       1465230523
                        1485130523
U107
       1465230523
                        1455130523
U108
       1465230523
                        1465230623
U109
       1465230523
                        1475130523
U110
        1465230523
                        1475130523
U111
        1465230523
                        1475130523
U112
        1465230523
                        1475130523
U113
        1465230523
                        1485130523
Time taken: 0.604 seconds, Fetched: 14 row(s)
```

```
hive> select * from station_geo_map;
OK
ST400
ST401
         AU
ST402
         AP
ST403
ST404
         E
ST405
ST406
         AU
ST407
         AP
ST408
        E
        E
ST409
ST410
ST411
ST412
         AP
ST413
Time taken: 0.447 seconds, Fetched: 14 row(s)
```

5. Data Formatting:

In this stage, we are merging the data coming from both web applications and mobile applications and create a common table for analysing purpose and create partitioned data based on batchid, since we are running this scripts for every 3 hours.

dataformatting.sh script:

```
DataFormatting.scala@
 1 import org.apache.spark.{SparkConf, SparkContext}
 2 import org.apache.spark.sql
 3
 4 object DataFormatting {
     def main(args: Array[String]): Unit = {
 5
        val conf = new SparkConf().setAppName("Data Formatting")
 6
 7
        val sc = new SparkContext(conf)
 8
        val sqlContext = new org.apache.spark.sql.hive.HiveContext(sc)
 9
        val batchId = args(0)
        val create hive table = """CREATE TABLE IF NOT EXISTS project.formatted input
10
11
12
                                          User_id STRING,
13
                                          Song id STRING,
14
                                          Artist id STRING,
15
                                          Timestamp STRING,
                                          Start_ts STRING,
16
17
                                          End_ts STRING,
                                          Geo cd STRING,
18
19
                                          Station id STRING,
20
                                          Song_end_type INT,
                                          Like INT,
21
                                          Dislike INT
22
23
24
                                          PARTITIONED BY
25
                                          (batchid INT)
26
                                          ROW FORMAT DELIMITED
27
                                          FIELDS TERMINATED BY ','
                                          ....
28
29
       val load mob data = s"""LOAD DATA LOCAL INPATH 'file:///home/acadgild/examples/music/data/mob/file.txt'
30
31
                                INTO TABLE project.formatted_input PARTITION (batchid='$batchId')"""
33
       val load_web_data = s"""INSERT INTO project.formatted_input
34
                                PARTITION (batchid='$batchId')
35
                                SELECT user_id,
36
                                song_id,
37
                                artist id.
                                unix_timestamp(timestamp,'yyyy-MM-dd HH:mm:ss') AS timestamp, unix_timestamp(start_ts,'yyyy-MM-dd HH:mm:ss') AS start_ts,
38
39
                                unix_timestamp(end_ts,'yyyy-MM-dd HH:mm:ss') AS end ts,
40
41
                                geo cd,
42
                                station id,
43
                                song_end_type,
44
                                like,
45
                                dislike
46
                                FROM web_data
47
48
49
51
        val xmlData = sqlContext.read.format("com.databricks.spark.xml").option("rowTag", "record").load("file:///home/acadgild/examples/music/data/web/file.xml"
52
        xmlData.createOrReplaceTempView("web_data")
53
54
55
        sglContext.sgl(load mob data)
        sqlContext.sql(load_web_data)
58
59
      catch{
      case e: Exception=>e.printStackTrace()
60
61 }
```

```
[acadgild@localhost music]$ cd MusicDataAnalysis
[acadgild@localhost MusicDataAnalysis]$ ls -ls
total 8
4 -rw-rw-r--. 1 acadgild acadgild 802 Dec 1 18:34 build.sbt
4 drwxrwxr-x. 3 acadgild acadgild 4096 Dec 1 18:34 src
```

Below is the command to create jar file in verbose mode:

```
[acadgild@localhost MusicDataAnalysis]$ sbt -v package
[process_args] java_version = '1.8'
# Executing command line:
java
-Xms1024m
-Xmx1024m
-XX:ReservedCodeCacheSize=128m
-XX:Rasmedsize=256m
-jar
/usr/share/sbt/bin/sbt-launch.jar
package

Getting org.scala-sbt sbt 1.0.4 (this may take some time)...
downloading https://repol.maven.org/maven2/org/scala-sbt/sbt/1.0.4/sbt-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#sbt;1.0.4!sbt.jar (910ms)
downloading https://repol.maven.org/maven2/org/scala-lang/scala-library/2.12.4/scala-library-2.12.4.jar ...
[SUCCESSFUL ] org.scala-lang#scala-library;2.12.4!scala-library.jar (22703ms)
downloading https://repol.maven.org/maven2/org/scala-sbt/Main 2.12/1.0.4/main 2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#main 2.12;1.0.4!main 2.12-jar (6809ms)
downloading https://repol.maven.org/maven2/org/scala-sbt/logic 2.12.jar (032ms)
downloading https://repol.maven.org/maven2/org/scala-sbt/actions 2.12/1.0.4/actions 2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#logic 2.12;1.0.4!logic 2.12.jar (1258ms)
downloading https://repol.maven.org/maven2/org/scala-sbt/main-settings_2.12/1.0.4/main-settings_2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#actions 2.12;1.0.4!main-settings_2.12.jar (2525ms)
```

Below is the location of Jar file which gets created under /MusicDataAnalysis/target/scala-2.11:

```
[acadgild@localhost MusicDataAnalysis] $ ls -ls
total 16

4 -rw-rw-r--. 1 acadgild acadgild 802 Dec 1 18:34 build.sbt

4 drwxrwxr-x. 3 acadgild acadgild 4096 Dec 1 18:52 project

4 drwxrwxr-x. 4 acadgild acadgild 4096 Dec 1 18:58 target

You have new mail in /var/spool/mail/acadgild
[acadgild@localhost MusicDataAnalysis] $ cd target
[acadgild@localhost target] $ ls -ls
total 8

4 drwxrwxr-x. 4 acadgild acadgild 4096 Dec 1 19:12 scala-2.11
[4 drwxrwxr-x. 4 acadgild acadgild 4096 Dec 1 18:53 streams
[acadgild@localhost target] $ cd scala-2.11
[acadgild@localhost target] $ ls -ls
total 16

4 drwxrwxr-x. 2 acadgild acadgild 4096 Dec 1 19:12 classes

8 -rw-rw-r--. 1 acadgild acadgild 8183 Dec 1 19:12 musicdataanalysis_2.11-1.0.jar

4 drwxrwxr-x. 5 acadgild acadgild 4096 Dec 1 19:10 resolution-cache
```

```
[acadgild@localhost MusicDataAnalysis]$ ls -ls
total 16
4 -rw-rw-r--. 1 acadgild acadgild 802 Dec
                                                   1 18:34 build.sbt
4 drwxrwxr-x. 3 acadgild acadgild 4096 Dec 1 18:52 project 4 drwxrwxr-x. 3 acadgild acadgild 4096 Dec 1 18:34 src 4 drwxrwxr-x. 4 acadgild acadgild 4096 Dec 1 18:58 target
[acadgild@localhost MusicDataAnalysis]$ cd src
[acadgild@localhost src]$ ls -ls
total 4
4 drwxrwxr-x. 3 acadgild acadgild 4096 Dec 1 18:34 main
[acadgild@localhost src]$ cd main
[acadgild@localhost main]$ ls -ls
total 4
4 drwxrwxr-x. 2 acadgild acadgild 4096 Dec 1 18:40 scala
[acadgild@localhost main]$ cd scala
[acadgild@localhost scala]$ ls -ls
total 20
8 -rw-rw-r--. 1 acadgild acadgild 4814 Dec
                                                  1 18:34 DataAnalysis.scala
4 -rw-rw-r--. 1 acadgild acadgild 3264 Dec 1 18:34 DataEnrichment.scala
4 -rw-rw-r--. 1 acadgild acadgild 2620 Dec 1 18:40 DataFormatting.scala
```

We are executing master script which internally calls **dataformatting.sh** which performs data formatting:

```
Perpairing to execute python scripts to generate data...

Data Generated Successfully |

Perpairing to execute python scripts to generate data...

Data Generated Successfully |

13921 ResourceHanager |
13921 ResourceHanager |
14931 JOBHISTORYSETYPE |
14031 Modelmanger |
14791 JOBHISTORYSETYPE |
14031 Modelmanger |
14791 JOBHISTORYSETYPE |
14032 SecondaryNameNode |
13092 Statabode |
13092 Sta
```

```
18/12/01 20:17:20 INFO metastore.HiveMetaStore: 0: get_database: project
18/12/01 20:17:20 INFO HiveMetaStore.audit: ugi=acadgld je_munknown.ip-addr cmd=get_database: project
18/12/01 20:17:20 INFO metastore.HiveMetaStore.audit: ugi=acadgld je_munknown.ip-addr cmd=get_database: project
18/12/01 20:17:20 INFO metastore.HiveMetaStore.audit: ugi=acadgld je_munknown.ip-addr cmd=get_table: db=project tbl=formatted_input
18/12/01 20:17:20 INFO metastore.HiveMetaStore.audit: ugi=acadgld je_munknown.ip-addr cmd=get_table: db=project tbl=formatted_input
18/12/01 20:17:20 INFO parser.CatalystSqlParser: Parsing command: int je_munknown.ip-addr cmd=get_table: db=project tbl=formatted_input
18/12/01 20:17:20 INFO parser.CatalystSqlParser: Parsing command: string
18/12/01 20:17:21 INFO parser.CatalystSqlParser: Parsing command: int
```

Below hive table **formatted_input** gets created which contains all data which gets merged from web and mobile applications (file.txt and file.xml):

```
nive> show tables;
formatted input
song_artist_map
station_geo_map
subscribed_users
rime taken: 0.221 seconds, Fetched: 4 row(s)
nive> select * from formatted_input;
                                                      1475130523
1465130523
                                                                            1465230523
1485130523
U120
          S203
                     A302
                                1495130523
                                                                                                            ST410
                                                                                                  AU
J106
          S203
                     A303
                                1495130523
                                                                                                             ST403
                                                      1485130523
                                                                            1475130523
                                                                                                             ST403
          S200
S202
                     A301
A305
                                1475130523
1475130523
                                                                            1485130523
1465130523
U108
                                                      1485130523
                                                                                                             ST410
J115
                                                      1475130523
                                                                                                             ST403
                     A304
                                1495130523
                                                      1485130523
                                                                            1475130523
                                                                                                             ST404
U101
U105
                                                                                                 AU
          S202
S208
                                1495130523
1465230523
                                                      1475130523
1465230523
                                                                            1485130523
1475130523
                                                                                                            ST406
ST400
                     A300
                     A301
                                1465230523
1465130523
1495130523
                                                      1465130523
1465130523
1475130523
J101
          5201
                     A302
                                                                            1475130523
                                                                                                             ST412
U112
U110
          S203
S209
                                                                            1475130523
1475130523
                                                                                                            ST406
                     A303
                                                                                                             ST406
                     A300
A301
A301
U100
          5207
                                1475130523
                                                      1485130523
                                                                            1485130523
                                                                                                             ST413
          S202
S203
                                1465130523
                                                      1475130523
1485130523
U103
                                                                            1485130523
                                                                                                             ST404
 109
                                1465130523
                                                                            1485130523
                                                                                                             ST415
U102
          5204
                     A301
                                1465230523
                                                      1485130523
                                                                            1475130523
                                                                                                             ST411
                                                      1465230523
U111
U107
          5200
                     A303
                                1495130523
                                                                            1465230523
                                                                                                             ST404
                     A301
                                1465130523
                                                      1475130523
                                                                            1465230523
                                                                                                             ST409
U114
          S210
                     A302
                                1465130523
                                                      1465230523
                                                                            1475130523
                                                                                                  A
AP
                                                                                                             ST409
U109
                     A301
                                1465230523
                                                      1485130523
                                                                            1485130523
          S200
                                                                                                             ST407
 110
          S200
                     A300
                                 1465230523
                                                      1485130523
                                                                            1475130523
                                                                                                             ST404
                                                                                                             ST407
ST415
U105
U100
          S205
                     A300
A304
                                1465490556
                                                      1462863262
                                                                            1462863262
1465490556
          S205
                                1468094889
                                                      1468094889
                                                                                                  AU
J100
          S203
                     A302
                                 1462863262
                                                       1468094889
                                                                            1465490556
                                                                                                             ST403
U119
          S202
                     A304
A305
                                1462863262
                                                      1465490556
                                                                            1462863262
                                                                                                  A
AP
                                                                                                             ST408
U114
          S210
                                                                            1465490556
                                                                                                             ST409
                                1494297562
                                                      1468094889
          S202
S204
                     A304
                                 1462863262
                                                                            1465490556
                                                       1462863262
                                                                                                             ST415
                     A300
                                1468094889
                                                      1494297562
                                                                            1494297562
                                                                                                             ST403
```

In the above screenshot we can see the formatted input data with some null values in **user_id, aritist_id** and **geo_cd** columns which we will fill the enrichment script based on rules of enrichment for **artist_id** and **geo_cd** only.

6. Data Enrichment and Cleaning:

In this phase we will enrich the data coming from web and mobile applications using the lookup table stored in HBase and divide the records based on the enrichment rules into 'pass' and 'fail' records.

- 1. If any of like or dislike is NULL or absent, consider it as 0.
- 2. If fields like Geo_cd and Artist_id are NULL or absent, consult the lookup tables for fields

 Station id and Song id respectively to get the values of Geo cd and Artist id.
- 3. If corresponding lookup entry is not found, consider that record to be invalid

So based on the enrichment rules we will fill the null geo_cd and artist_id values with the help of corresponding lookup values in song-artist-map and station-geo-map tables in Hive-Hbase tables.

```
******************
#!/bin/bash
batchid=`cat /home/acadgild/examples/music/logs/current-batch.txt`
LOGFILE=/home/acadgild/examples/music/logs/log batch $batchid
VALIDDIR=/home/acadgild/examples/music/processed dir/valid/batch $bat
INVALIDDIR=/home/acadgild/examples/music/processed dir/invalid/batch
$batch id
echo "Running script for data enrichment and filtering..." >> $LOGFILE
spark-submit --class DataEnrichment \
--master local[2] \
--jars /home/acadqild/install/hive/apache-hive-2.3.2-bin/lib/hive-
hbase-
handler-2.3.2.jar,/home/acadgild/install/hive/apache-hive-2.3.2-
bin/lib/hbase-client-
1.1.1.jar,/home/acadgild/install/hive/apache-hive-
2.3.2-bin/lib/hbase-common-
1.1.1.jar,/home/acadgild/install/hive/apache-
hive-2.3.2-bin/lib/hbase-hadoop-compat-
1.1.1.jar,/home/acadgild/install/hive/apache-hive-2.3.2-bin/lib/hbase-
server-1.1.1.jar,/home/acadgild/install/hive/apache-hive-2.3.2-
bin/lib/hbase-protocol-
```

1.1.1.jar,/home/acadgild/install/hive/apache-hive-

```
2.3.2-bin/lib/zookeeper-3.4.6.jar,/home/acadgild/install/hive/apache-
hive-
2.3.2-bin/lib/guava-14.0.1.jar,/home/acadgild/install/hive/apache-
hive-
2.3.2-bin/lib/htrace-core-3.1.0-incubating.jar \
/home/acadgild/examples/music/MusicDataAnalysis/target/scala-
2.11/musicdataanalysis 2.11-1.0.jar $batchid
if [ ! -d "$VALIDDIR" ]
then
mkdir -p
"$VALIDDIR" fi
if [ ! -d "$INVALIDDIR" ]
then
mkdir -p "$INVALIDDIR"
fi
echo "Copying valid and invalid records in local file system..." >>  
$LOGFILE
hadoop fs -get
/user/hive/warehouse/project.db/enriched data/batchid=$batchid/status=
pass/
* $VALIDDIR
hadoop fs -get
/user/hive/warehouse/project.db/enriched data/batchid=$batchid/status=
fail/
* $INVALIDDIR
echo "Deleting older valid and invalid records from local file
system..."
>> $LOGFILE
find /home/acadgild/examples/music/processed dir/ -mtime +7 -exec rm
{} \;
```

*

We have executed data_enrichment.sh script by calling music_project_master.sh batch file as shown below:

```
[acadgild@localhost music]s [./music_project_master.sh
Preparing to execute python scripts to generate data...

Data Generated Successfully !
Starting the daemons....

15888 RunJar
4528 HMaster
4528 HMaster
3890 NodeMananger
3990 NodeMananger
3990 NameNode
5011 RunJar
3555 SecondaryNameNode
4714 JobhistoryServer
4635 HMegloinServer
4636 HMegloinServer
4640 HQuorumPeer
17007 Jps
3791 ResourceManager
All hadoop daemons started !
Upload the look up tables now in Hbase...
Done with data population in look up tables !
Lets do some data formatting now....
data formatting complete !
Creating hive tables on top of hbase tables for data enrichment and filtering...
Hive table with Hbase Manoirno Comolete !
Let us do data enrichment as per the requirement...
Let with data formatting complete !
Creating hive tables on top of hbase tables for data enrichment and filtering ...
Hive table with Hbase Manoirno Comolete !
Let us do data enrichment as per the requirement...
Let us do data enrichment as per the requirement...
Let us do data enrichment as per the requirement...

1871/2/02 15:25:32 MARN util. MativecodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable ...
1871/2/02 15:25:33 MARN util. MativecodeLoader: Unable to load native-hadoop library for your platform... using 192.168.0.100 in stead (on interface ethls)
1871/2/02 15:25:33 MARN util. MativecodeLoader: Changing view acts to: acadgild
1871/2/02 15:25:33 No spark SecurityManager: changing wearst to: acadgild
1871/2/02 15:25:33 No spark SecurityManager: Changing modify acts to: acadgild
1871/2/02 15:25:33 No spark SecurityManager: Changing modify acts groups to:
1871/2/02 15:25:33 No spark SecurityManager: Changing modify acts groups to:
1871/2/02 15:25:33 No spark SecurityManager: Changing modify acts groups to:
1871/2/02 15:25:33 No spark SecurityManager: Changing year acts to: acadgild
1871/2/02 15:25:33 No spark SecurityManager: Changing year acts groups to:
```

```
18/12/02 15:28:44 INFO HiveMetaStore.audit: ugi=acadgild ip=unknown.ip-addr cmd=get_database: project
18/12/02 15:28:44 INFO metastore.kiveMetaStore: 0: get_table : db=project tbl=enriched_data
18/12/02 15:28:44 INFO metastore.audit: ugi=acadgild ip=unknown.ip-addr cmd=get_table : db=project tbl=enriched_data
18/12/02 15:28:44 INFO miscastore.audit: ugi=acadgild ip=unknown.ip-addr cmd=get_table : db=project tbl=enriched_data
18/12/03 15:28:44 INFO miscastore.audit: ugi=acadgild ip=unknown.ip-addr cmd=get_table : db=project tbl=enriched_data
18/12/02 15:28:44 INFO parser.clatalystSqlParser: Parsing command: int
18/12/02 15:28:44 INFO parser.clatalystSqlParser: Parsing command: string
18/12/02 15:28:45 INFO parser.clatalystSqlParser: Parsing command: string
18/12/02 15:28:45 INFO
```

Data enrichment tables:

```
hive> show databases;

OK
default
project
Time taken: 4.22 seconds, Fetched: 2 row(s)
hive> use project;

OK
Time taken: 0.116 seconds
hive> show tables;

OK
enriched_data
formatted_input
song_artist_map
station_geo_map
subscribed_users
Time taken: 0.251 seconds, Fetched: 5 row(s)
```

In the below screenshot, we have data for data enrichment table where we filled the null values of artist_id and geo_cd of formatted input with the help of lookup tables

hive> OK	select *	from en	riched_data;									
U111	S201	A301	1465490556	1494297562	1465490556	J	ST403	1	1	1	1	fail
U101	S201	A301	1465230523	1465130523	1475130523	AP	ST412	1	0	0	1	fail
U100	S207	A303	1475130523	1485130523	1485130523	J	ST413	1	1	1	1	fail
U103	5202	A302	1465130523	1475130523	1485130523	E	ST404	1	1	1	1	fail
U119	S202	A302	1462863262	1465490556	1462863262	E	ST408	3	1	1	1	fail
NULL	5202	A302	1462863262	1462863262	1465490556	NULL	ST415	Θ	1	1	1	fail
	S206	A302	1495130523	1485130523	1475130523	E	ST404	1	1	1	1	fail
U105	5208	A304	1465230523	1465230523	1475130523	A	ST400		1	1	1	fail
U114	S210	NULL	1465130523	1465230523	1475130523	E	ST409	0	0	1	1	fail
U114	S210	NULL	1494297562	1468094889	1465490556	E	ST409	2	1	Θ	1	fail
U108	S205	A301	1462863262	1468094889	1465490556	A	ST410	1	1	1	1	fail
U105	S205	A301	1465490556	1462863262	1462863262	AP	ST407	0	1	1	1	fail
U100	S205	A301	1468094889	1468094889	1465490556	NULL	ST415	2	Θ	1	1	fail
U110	S200	A300	1465230523	1485130523	1475130523	E	ST404	1	1	1	1	fail
U113	5203	A303	1465490556	1465490556	1468094889	AP	ST407	0	0	Θ	1	fail
U109	S203	A303	1465130523	1485130523	1485130523	NULL	ST415	1	1	0	1	fail
U114	S203	A303	1494297562	1462863262	1468094889	NULL	ST415	3	1	0	1	fail
U112	S203	A303	1465130523	1465130523	1475130523	AU	ST406	0	1	1	1	fail
U106	S201	A301	1468094889	1462863262	1462863262	J	ST403	2	Θ	1	1	pass
U106	S207	A303	1494297562	1494297562	1468094889	E	ST404		Θ	1	1	pass
U117	S202	A302	1462863262	1465490556	1465490556	E	ST404	0	1	0	1	pass
U115	5202	A302	1475130523	1475130523	1465130523	J	ST403	2	Θ	Θ	1	pass
U101	S202	A302	1495130523	1475130523	1485130523	AU	ST406		0	1	1	pass
U102	5204	A304	1494297562	1462863262	1465490556	E	ST414		1	0	1	pass
U119	5204	A304	1475130523	1485130523	1475130523	J	ST403	0	0	1	1	pass
U109	S204	A304	1468094889	1494297562	1494297562	J	ST403		Θ	1	1	pass
U103	S204	A304	1462863262	1465490556	1465490556	Α	ST410	3	0	1	1	pass
U102	S204	A304	1465230523	1485130523	1475130523	A	ST411	0	Θ	0	1	pass
U104	5209	A305	1465490556	1462863262	1494297562	AP	ST407	0	0	1	1	pass
U110	S209	A305	1495130523	1475130523	1475130523	AU	ST406	0	1	0	1	pass
U116	S206	A302	1465490556	1462863262	1468094889	E	ST409		1	Θ	1	pass
U118	S206	A302	1465490556	1465490556	1462863262	A	ST411	1	0	1	1	pass
U107	S205	A301	1465130523	1475130523	1465230523	E	ST409	1	1	0	1	pass
U104	S205	A301	1462863262	1468094889	1468094889	Е	ST409	2	Θ	0	1	pass

At the end, script will automatically divide the records based on status **pass & fail** and dump the result into **processed_dir** folder with **valid** and **invalid** folders as shown below:

Enrichment phase is executed successfully by applying all the rules of enrichment.

7. Data Analysis:

In this stage, we will do analysis on enriched data using Spark SQL and run the program using **Spark-Submit** command.

Data_analysis.sh script file:

```
*********

#!/bin/bash

batchid=`cat /home/acadgild/examples/music/logs/current-batch.txt`
LOGFILE=/home/acadgild/examples/music/logs/log_batch_$batchid

echo "Running script for data analysis..." >> $LOGFILE

spark-submit --class DataAnalysis --master local[2] \
--jars /home/acadgild/install/hive/apache-hive-2.3.2-bin/lib/hive-hbase-handler-2.3.2.jar,/home/acadgild/install/hive/apache-hive-2.3.2-bin/lib/hbase-client-
```

```
1.1.1.jar,/home/acadgild/install/hive/apache-hive-2.3.2-
bin/lib/hbase-common-1.1.1.jar,/home/acadgild/install/hive/apache-
hive-2.3.2-bin/lib/hbase-hadoop-compat-
1.1.1.jar,/home/acadgild/install/hive/apache-hive-2.3.2-
bin/lib/hbase-server-1.1.1.jar,/home/acadgild/install/hive/apache-
hive-2.3.2-bin/lib/hbase-protocol-
1.1.1.jar,/home/acadgild/install/hive/apache-hive-2.3.2-
bin/lib/zookeeper-3.4.6.jar,/home/acadgild/install/hive/apache-hive-
2.3.2-bin/lib/guava-14.0.1.jar,/home/acadgild/install/hive/apache-
hive-2.3.2-bin/lib/htrace-core-3.1.0-incubating.jar \
/home/acadgild/examples/music/MusicDataAnalysis/target/scala-
2.11/musicdataanalysis 2.11-1.0.jar $batchid
sh /home/acadgild/examples/music/data export.sh
echo "Incrementing batchid..." >> $LOGFILE
batchid=`expr $batchid + 1`
echo -n $batchid > /home/acadgild/examples/music/logs/current-
batch.txt
DataAnalysis.scala Program:
DataAnalysis.scala program:
import org.apache.spark.{SparkConf, SparkContext}
import org.apache.spark.sql
object DataAnalysis {
def main(args: Array[String]): Unit = {
val conf = new SparkConf().setAppName("Data Analysis")
val sc = new SparkContext(conf)
val sqlContext = new
org.apache.spark.sql.hive.HiveContext(sc)
val batchId = args(0)
```

// Problem 1 :Determine top 10 station_id(s) where
maximum number of songs were played, which were liked

by unique users.

```
val create top 10 stations = """CREATE TABLE IF NOT
EXISTS
top_10_stations
station id STRING, total distinct songs played INT,
distinct user count INT
)
PARTITIONED BY (batchid INT) ROW FORMAT DELIMITED
FIELDS TERMINATED BY ',' STORED AS TEXTFILE"""
val load_top_10_stations = s"""INSERT OVERWRITE TABLE
top 10 stations
PARTITION (batchid='$batchid') SELECT
station id,
COUNT (DISTINCT song id) AS total distinct songs played,
COUNT (DISTINCT user id) AS distinct user count
FROM enriched data
WHERE status='pass'
AND batchid='$batchId'
AND like=1
GROUP BY station id
ORDER BY total distinct songs played DESC LIMIT 10"""
// Problem 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.
val create users behaviour = """CREATE TABLE IF NOT
EXISTS
users behaviour
(
user type STRING, duration INT
)
PARTITIONED BY (batchid INT) ROW FORMAT DELIMITED
```

```
FIELDS TERMINATED BY ','
STORED AS TEXTFILE"""
val load users_behaviour = s"""INSERT OVERWRITE TABLE
users behaviour
PARTITION (batchid='$batchId') SELECT
CASE WHEN (su.user id IS NULL OR CAST(ed.timestamp AS
DECIMAL (20,0))
> CAST(su.subscn end dt AS DECIMAL(20,0))) THEN
'UNSUBSCRIBED'
WHEN (su.user id IS NOT NULL AND CAST(ed.timestamp AS
DECIMAL(20,0))
<= CAST(su.subscn end dt AS DECIMAL(20,0))) THEN
'SUBSCRIBED' END AS user type,
SUM(ABS(CAST(ed.end ts AS DECIMAL(20,0))-
CAST(ed.start ts AS DECIMAL(20,0)))) AS duration
FROM enriched data ed
LEFT OUTER JOIN subscribed users su
ON ed.user id=su.user id
WHERE ed.status='pass'
AND ed.batchid='$batchId'
GROUP BY CASE WHEN (su.user id IS NULL OR
CAST(ed.timestamp AS
DECIMAL(20,0)) > CAST(su.subscn end dt AS
DECIMAL(20,0)) THEN
'UNSUBSCRIBED'
WHEN (su.user id IS NOT NULL AND CAST(ed.timestamp AS
DECIMAL (20,0))
<= CAST(su.subscn end dt AS DECIMAL(20,0))) THEN
'SUBSCRIBED' END"""
//Problem 3 : Determine top 10 connected artists.
Connected artists are those whose songs are most
listened by the unique users who follow them.
val create connected artists = """CREATE TABLE IF NOT
EXISTS
connected artists
(
```

artist id STRING, user count INT

```
)
PARTITIONED BY (batchid INT) ROW FORMAT DELIMITED
FIELDS TERMINATED BY ',' STORED AS TEXTFILE"""
val load connected artists = s"""INSERT OVERWRITE TABLE
connected artists
PARTITION (batchid='$batchid') SELECT
ua.artist id,
COUNT (DISTINCT ua.user id) AS user count
FROM (
SELECT user id, artist id FROM users artists
LATERAL VIEW explode (artists array) artists AS
artist id
) ua
INNER JOIN (
SELECT artist id, song id, user id
FROM enriched data
WHERE status='pass'
AND batchid='$batchId'
) ed
ON ua.artist id=ed.artist id AND ua.user id=ed.user id
GROUP BY ua.artist id
ORDER BY user count DESC LIMIT 10"""
//Problem 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.
val create top 10 royalty songs = """CREATE TABLE IF
NOT EXISTS
top 10 royalty songs
song id STRING, duration INT
PARTITIONED BY (batchid INT) ROW FORMAT DELIMITED
FIELDS TERMINATED BY ',' STORED AS TEXTFILE"""
```

```
val load top 10 royalty songs = s"""INSERT OVERWRITE
TABLE
top 10 royalty songs
PARTITION (batchid='$batchId') SELECT song id,
SUM(ABS(CAST(end ts AS DECIMAL(20,0))-CAST(start ts AS
DECIMAL(20,0)))) AS duration
FROM enriched data
WHERE status='pass'
AND batchid='$batchId'
AND (like=1 OR song end type=0) GROUP BY song id
ORDER BY duration DESC LIMIT 10"""
//Problem 5: Determine top 10 unsubscribed users who
listened to the songs for the longest duration.
val create top 10 unsubscribed users = """CREATE TABLE
IF NOT EXISTS
top 10 unsubscribed users
user id STRING, duration INT
PARTITIONED BY (batchid INT) ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE"""
val load top 10 unsubscribed users = s"""INSERT
OVERWRITE TABLE
top 10 unsubscribed users
PARTITION(batchid='$batchId') SELECT
ed.user id,
SUM(ABS(CAST(ed.end ts AS DECIMAL(20,0))-
CAST(ed.start ts AS DECIMAL(20,0)))) AS duration
FROM enriched data ed
LEFT OUTER JOIN subscribed users su
```

ON ed.user id=su.user id

```
WHERE ed.status='pass'
AND ed.batchid='$batchId'
AND (su.user id IS NULL OR (CAST(ed.timestamp AS
DECIMAL(20,0)) > CAST(su.subscn end dt AS
DECIMAL(20,0))))
GROUP BY ed.user id ORDER BY duration DESC LIMIT 10"""
try {
sqlContext.sql("SET hive.auto.convert.join=false")
sqlContext.sql("USE project")
sqlContext.sql(create top 10 stations)
sqlContext.sql(load top 10 stations)
sqlContext.sql(create_users_behaviour)
sqlContext.sql(load users behaviour)
sqlContext.sql(create connected artists)
sqlContext.sql(load connected artists)
sqlContext.sql(create top 10 royalty songs)
sqlContext.sql(load top 10 royalty songs)
sqlContext.sql(create_top_10_unsubscribed users)
sqlContext.sql(load top 10 unsubscribed users)
}
catch{
case e: Exception=>e.printStackTrace()
}
}
}
*******************
```

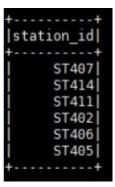
We are executing **Data analysis.sh** script by running **music project master.sh** script file.

```
18/12/99 15:45:58 INFO parser.CatalystSqlParser: Parsing command: string
18/12/99 15:45:58 INFO metastore.HiveMetaStore: 0: get_table: db=project_tbl=top_10_unsubscribed_users
18/12/99 15:45:58 INFO parser.CatalystSqlParser: Parsing command: int
18/12/99 15:45:59 INFO datasources.SqLHadopMapReduceCommitProtocol: Using output committer class org.apache.hadoop.mapreduce.lib.output.
18/12/99 15:45:59 INFO datasources.SqLHadopMapReduceCommitProtocol: Using output committer class org.apache.hadoop.mapreduce.lib.output.
18/12/99 15:45:59 INFO gagregate.HashAggregateExec: spark.sql.codegen.aggregate.map.twolevel.enable is set to true, but current version of codegened fast hashmap does not support this aggregate.
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 439
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 439
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 433
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 437
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 427
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 437
18/12/99 15:45:59 INFO spark.contextCleaner: cleaned accumulator 434
18/12/99
```

```
Warning: /home/acadgild/install/sqoop/sqoop-1.4.6.bin_ hadoop-2.0.4-alpha/../hcatalog does not exist! HCatalog jobs will fail.
Please set SHCAT_HOME to the root of your HCatalog installation.
Warning: /home/acadgild/install/sqoop/sqoop-1.4.6.bin_ hadoop-2.0.4-alpha/../accumulo does not exist! Accumulo imports will fail.
Please set SACCUMULO HOME to the root of your Accumulo installation.
18/12/09 16:21:09 NFG sqoop.Sqoop: Running Sqoop version: 1.4.6
18/12/09 16:21:09 NFG sqoop.Sqoop: Running Sqoop version: 1.4.6
18/12/09 16:21:10 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
18/12/09 16:21:10 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
18/12/09 16:21:10 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
18/12/09 16:21:10 INFO manager.Square.SSL connection without server's identity verification is not recommended. According to My SQL 5.5.454, 5.6.264 and 5.7.64 requirements SSL connection must be established by default if explicit option isn't set. For compliance with existing applications not using SSL the verifyserverCertificate property is set to 'false'. You need either to explicitly disable SSL by setting useSSL-false, or set useSSL-frue and provide truststore for server certificate verification.
18/12/09 16:21:19 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM 'top_10 stations' AS t LIMIT 1
18/12/09 16:21:19 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM 'top_10 stations' AS t LIMIT 1
18/12/09 16:21:19 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM 'top_10 stations' AS t LIMIT 1
18/12/09 16:21:19 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM 'top_10 stations' AS t LIMIT 1
18/12/09 16:21:31 INFO orm.CompilationManager: HADOOP MAPRED HOME is /home/acadgild/install/hadoop/hadoop-2.6.5
Note: Recompile with -Xlint:deprecation for details.
18/12/09 16:21:32 INFO maperduce.ExportJobBase: Beginning export of top_10 stations
SLF41: Seen http://www.slf41.org/codes
```

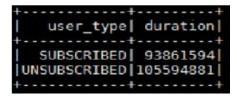
Solution1:

Determine top 10 station_id(s) where maximum number of songs were played, which were liked by unique users.



Solution2:

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.



Solution3:

Determine top 10 connected artists. Connected artists are those whose songs are most listened by the unique users who follow them.



Solution4:

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

```
| song_id
| s208
| s207
| s206
| s209
| s209
| s200
| s204
| s202
| s205
```

Solution5:

Determine top 10 unsubscribed users who listened to the songs for the longest duration.

```
| user_id|
| user_id|
| ull7|
| ull8|
| ull0|
| ul20|
| ul15|
| ul07|
| ul08|
| ul09|
| ul06|
| ul00|
```

We could see below that all tables have also been created in the Hive :

```
hive> use project;

OK

Time taken: 0.098 seconds
hive> show tables;

OK

connected artists
enriched_data
formatted_input
song_artist_map
station_geo_map
subscribed users
top_10_royalty_songs
top_16_stations
top_10_unsubscribed_users
users artists
users_behaviour

Time taken: 0.407 seconds, Fetched: 11 row(s)
hive>
```

The data analysis result is shown in the Hive tables below in the screen shot:

Below is the output of **top_10_stations** table: Below is the output of **users_behaviour** table: Below is the output of **connected artists** table:

```
Nive> Select * From top_10_stations;

OK

top_10_stations.station_id top_10_stations.total_distinct_songs_played top_10_stations.distinct_user_count top_10_stations.batchid

ST407 2 3 1

ST414 1 1 1

ST411 1 1 1

ST402 1 2 1

ST406 1 1 1

ST405 1 1 1

Time taken: 0.336 seconds, Fetched: 6 row(s)
```

```
hive> Select * From users_behaviour;

OK

users_behaviour.user_type users_behaviour.duration users_behaviour.batchid

SUBSCRIBED 93861594 1

UNSUBSCRIBED 105594881 1

Time taken: 0.274 seconds, Fetched: 2 row(s)
```

```
hive> Select * From connected_artists;

OK
connected_artists.artist_id connected_artists.user_count connected_artists.batchid

A303 2 1

A302 2 1

A300 1 1

Time taken: 0.225 seconds, Fetched: 3 row(s)
```

```
hive> Select * From top_10_royalty_songs;
OK
   _10_royalty_songs.song_id
                                 top_10_royalty_songs.duration
                                                                   top_10_royalty_songs.batchid
top_
S208
        22627294
S207
        20000000
5206
        19900000
5209
        15254588
5200
        9900000 1
        2604333
5204
5202
        100000
5205
Time taken: 0.237 seconds, Fetched: 8 row(s)
```

Now we need to export all the data to the MYSQL using sqoop, by executing data export.sh script file:

By using **data_export.sh** script file, we are going to export the data from the hive tables into mysql using Sqoop export.

```
data_export.sh
1 #!/bin/bash
3 #This script is not working.
 4 #Either change table to text or use STRING as type of partitioned column
 6 batchid='cat /home/acadgild/examples/music/logs/current-batch.txt
 7 LOGFILE=/home/acadgild/examples/music/logs/log batch $batchid
9 echo "Creating mysgl tables if not present...AnkithTest" >> $LOGFILE
11 mysql -u "root" "-pRoot@123" < /home/acadgild/examples/music/create_schema.sql
13 echo "Running sqoop job for data export...AnkithTest" >> $LOGFILE
15 sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table top 10 stations --export-dir /user/hive/warehouse/project.db/
/top_10_stations/batchid=$batchid --input-fields-terminated-by ',' -m 1
17 sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table users_behaviour --export-dir /user/hive/warehouse/project.db/
//users_behaviour/batchid=Sbatchid --input-fields-terminated-by ',' -m 1
19 sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table connected_artists --export-dir /user/hive/warehouse/project.db/
/connected artists/batchid=$batchid --input-fields-terminated-by ',' -m 1
21 sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table top 10_royalty_songs --export-dir /user/hive/warehouse/project.db/
/top_10_royalty_songs/batchid=$batchid --input-fields-terminated-by ',' -m 1
23 sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table top_10_unsubscribed_users --export-dir /user/hive/warehouse/
/project.db/top_10_unsubscribed_users/batchid=$batchid --input-fields-terminated-by ',' -m 1
  create_schema.sql
   1 CREATE DATABASE IF NOT EXISTS project;
   2
   3 USE project;
   5 CREATE TABLE IF NOT EXISTS top_10_stations
   7 station_id VARCHAR(50),
8 total_distinct_songs_ploatinct_user_count INT
                                                   played INT,
 10 ) 2
  11
 12 CREATE TABLE IF NOT EXISTS users_behaviour
 13 (
                type VARCHAR (50),
 14 user
  15 duration BIGINT
 16 ) 2
 17
  18 CREATE TABLE IF NOT EXISTS connected_artists
 19
 20 artist_id VARCHAR (50) ,
      user_count INT
 22 ) ;
 23
24 CREATE TABLE IF NOT EXISTS top_10_royalty_songs
  25
  26 song_id VARCHAR(50),
27 duration BIGINT
 26 song_id VARCHAR(50),
27 duration BIGINT
28 );
29
30 CREATE TABLE IF NOT EXISTS top_10_unsubscribed_users
31 (
32 user_id VARCHAR(50),
33 duration BIGINT
34 );
35 |
 36 commit:
```

Below we could see that data exported successfully into the MYSQL Database for all the 5 queries:

The sqoop export command exported the tables from the hive and it stored in the Mysql. The below screen shot show the successful Sqoop export from hive to mysql. The data stored in the Mysql is shown in below screenshots:

```
Warning: /home/acadgild/install/sqoop/sqoop-1.4.6.bin_hadoop-2.0.4-alpha/../hcatalog does not exist! HCatalog jobs will fail.
Please set sHCAT_HOM€ to the root of your HCatalog installation.
Warning: /home/acadgild/install/sqoop/sqoop-1.4.6.bin_hadoop-2.0.4-alpha/../accumulo does not exist! Accumulo imports will fail.
Please set sACCUMNLO hOME to the root of your Accumulo installation.
18/12/99 16:21:09 IMF0 sqoop.Sqoop: Running Sqoop version: 1.4.6
18/12/99 16:21:09 IMF0 sqoop.Sqoop: Running Sqoop version: 1.4.6
18/12/99 16:21:10 IMF0 manager.MySQI Manager: Preparing to use a MySQI streaming resultset.
18/12/99 16:21:11 IMF0 tool.CodeGenTool: Beginning code generation
SUN Doc 09 16:21:12 IST 2018 WARN: Establishing SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server's identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection without server identity verification is not recommended. According to My
SQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL 5.40+ requirements of the server requirements of the server requirements of the server requirements of t
```

The **project** database had been exported from hive (HDFS) and the below screen shot shows all tables:

```
mysql> use project;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
mysql> show tables;

| Tables_in_project |
| connected_artists |
| song_duration |
| top_10_royalty_songs |
| top_10_stations |
| top_10_unsubscribed_users |
| top_10_unsubscribed_users |
| top_10_unsubscribed_users |
```

Output from **top_10_stations** table in mysql is shown below:

Output from users_behaviour table in mysgl is shown below:

```
mysql> Select * From users_behaviour;
| user_type | duration |
| SUBSCRIBED | 93861594 |
| UNSUBSCRIBED | 105594881 |
| the state of th
```

Output from **connected_artists** table in mysql is shown below:

Output from **top_10_royalty_songs** table in mysql is shown below:

```
top_10_royalty_songs;
song_id
           duration
           22627294
S208
S207
           20000000
           19900000
S206
           15254588
S209
            9900000
S200
S204
            2604333
             100000
S202
S205
                   Θ
rows in set (0.00 sec)
```

Output from top_10_unsubscribed_users table in mysql is shown below:

```
mysql> Select * From top_10_unsubscribed_users;
  user_id | duration
  U117
           20000000
  U118
           20000000
 U110
           2000000
 U120
           12627294
           12527294
 U115
 U107
           10000000
 U108
             5231627
 U109
            2604333
 U106
             2604333
                   0
  U100
10 rows in set (0.01 sec)
```

8. Job Scheduling

Now after exporting data into MySQL, **batchid** will be incremented to additional 1 means one batch of data operations is successfully completed and new batch of data will be loaded for the analysis after every 3 hours.

We can check logs to track the behaviour of the operations we have done on the data and overcome failures (if any) we could see the **batchid** gets incremented by 1 in **current-batch.txt**

```
[acadgild@localhost logs]$ pwd
/home/acadgild/examples/music/logs
[acadgild@localhost logs]$ ls -ls
total 52
4 -rwxrwxr-x. 1 acadgild acadgild 2 Dec 9 17:18 current-batch.txt
4 -rw-rw-r--. 1 acadgild acadgild 522 Dec 9 16:21 log batch 1
```

```
[acadgild@localhost logs]$ cat current-batch.txt
```

Finally Jar file gets created as highlighted below: