

Music Data Analysis using Spark, Hive, HBase , Sqoop

- Prachi Mohite

Section – 1 - Project Overview

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

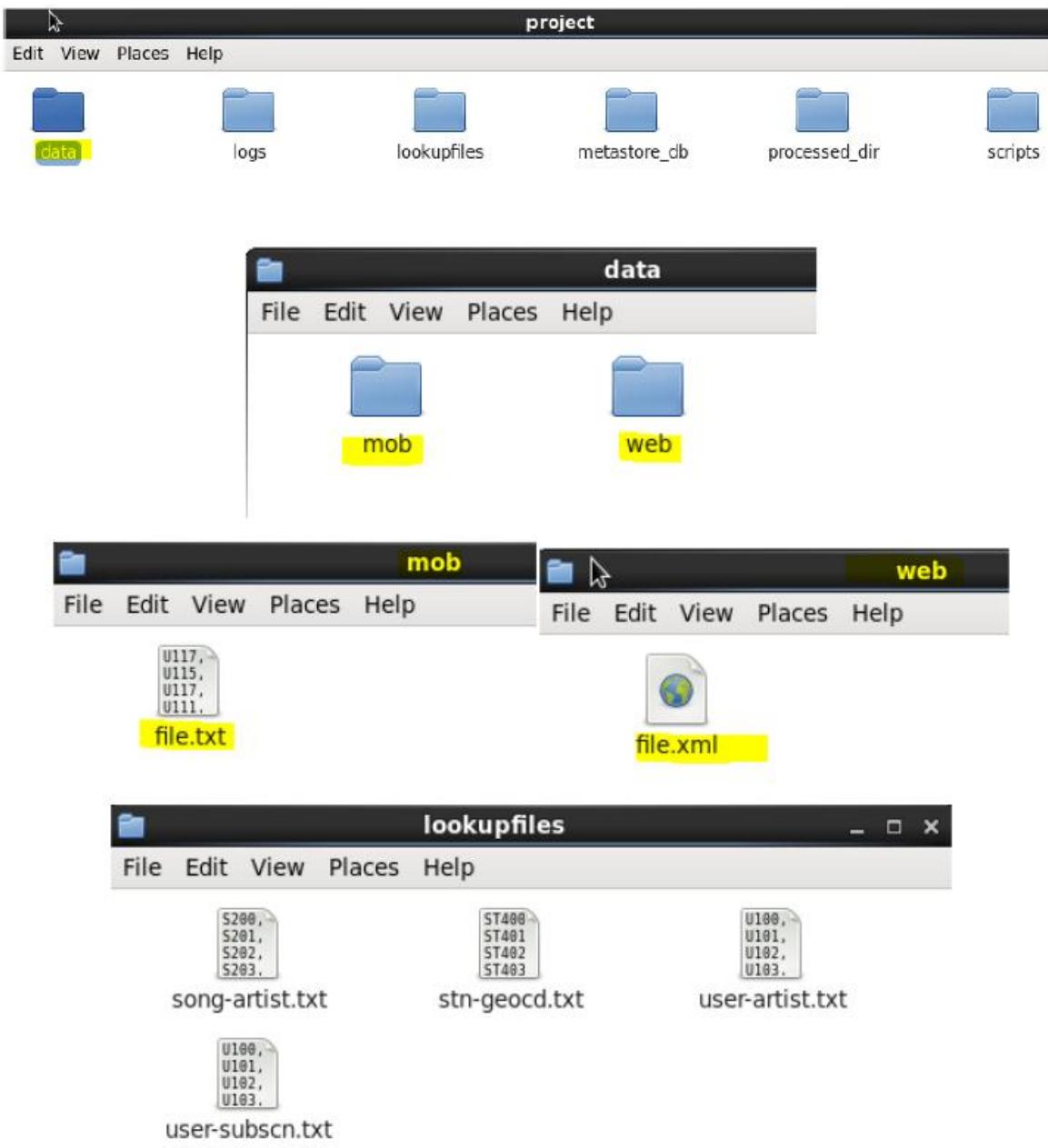
Table Name	Description
Station_Geo_Map	Contains mapping of a geo_cd with station_id
Subscribed_Users	Contains user_id, subscription_start_date and subscription_end_date. Contains details only for subscribed users
Song_Artist_Map	Contains mapping of song_id with artist_id alongwith royalty associated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

1.3 DATASET

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 present in lookup directory should be used in HBase.

Below is the link for same.

https://drive.google.com/drive/folders/0B_P3pWagdIrrMjGVINsSUEtbG8?usp=sharing



1.4 Data Enrichment

Rules for 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 (SHOULD BE IMPLEMENTED IN SPARK)

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

1.6 Challenges and Optimizations:

1. LookUp tables are in NoSQL databases. Integrate them with the actual data flow.
2. Try to make joins as less expensive as possible.
3. Data Cleaning, Validation, Enrichment, Analysis and Post Analysis have to be automated. Try using schedulers.
4. Appropriate logs have to maintain to track the behavior and overcome failures in the pipeline.

1.7 Flow of operations

A schematic flow of operations is shown below,

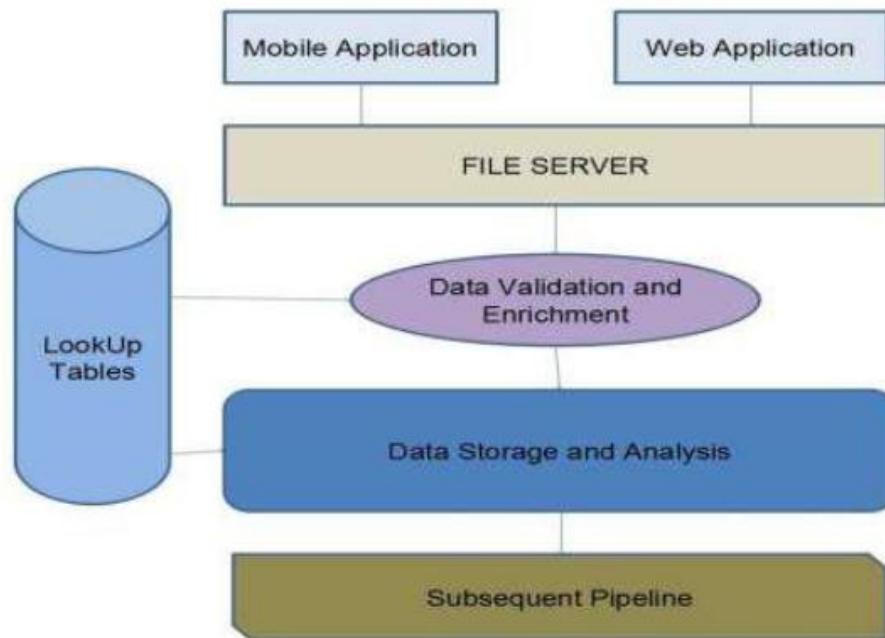


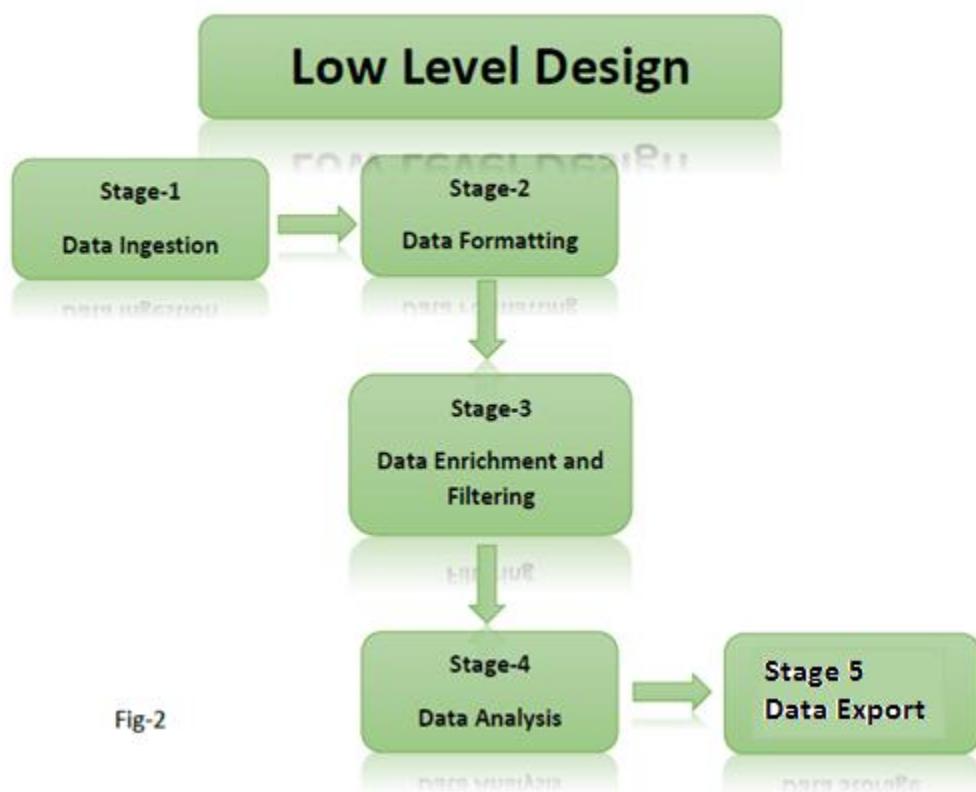
Fig-1

In the following sections, we are going to see the Music Data Analysis as per the above rules.

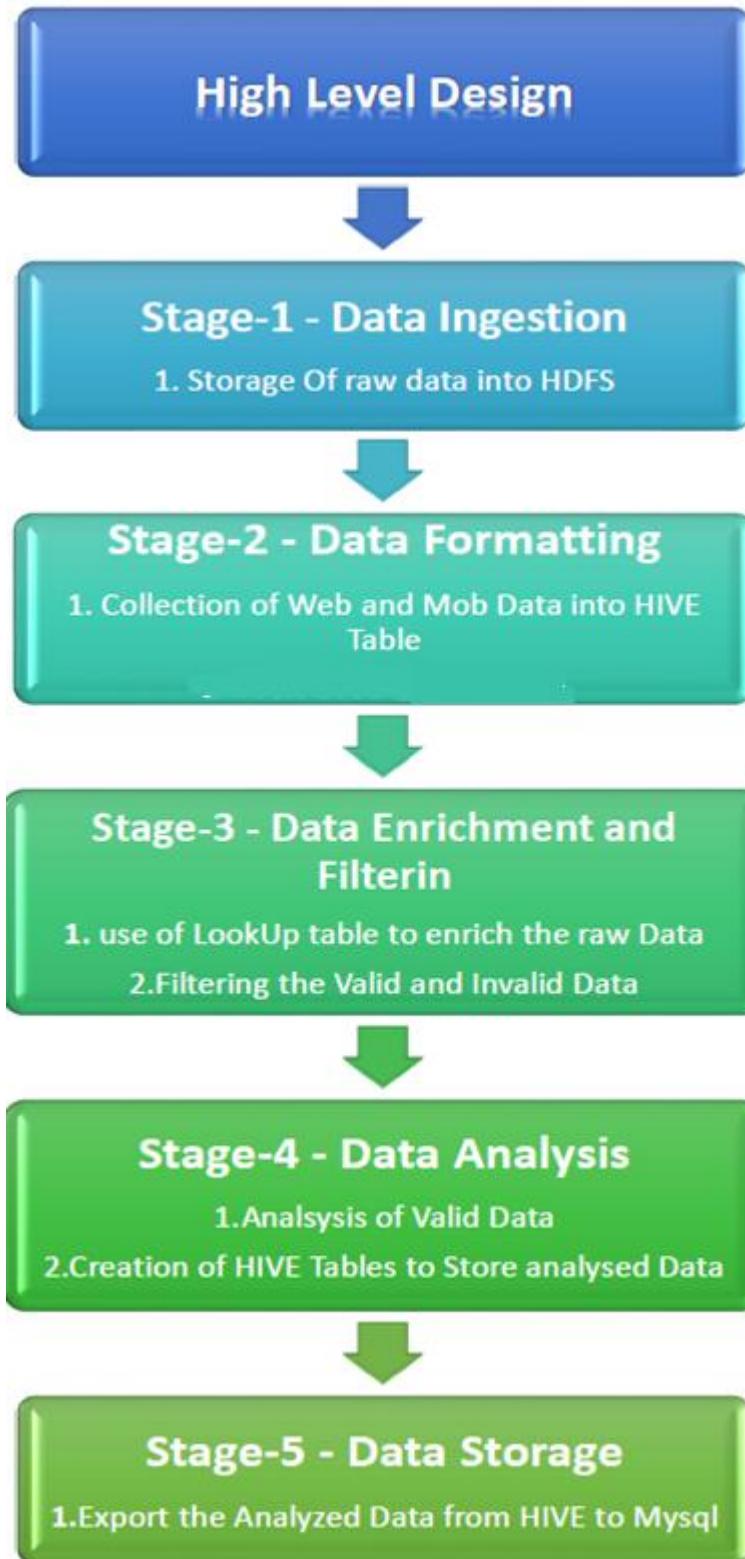
Section -2 – Design of the Project

2.1 Low Level Design

The following flowchart shows the Low Level design of this project,



2.2 High Level Design



Section-3-Hadoop Eco-System Implementation

1. We have created a batch file “**start-daemon.sh**” which starts the daemons such as **hive**, **hbase**, **Mysql** and rest of the all **hadoop** daemons.

Batch file script,

```
#!/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" ]
then
  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

# To start all Hadoop Daemons
start-all.sh
# To start hbase Daemons
start-hbase.sh
# To start job history server|
mr-jobhistory-daemon.sh start historyserver

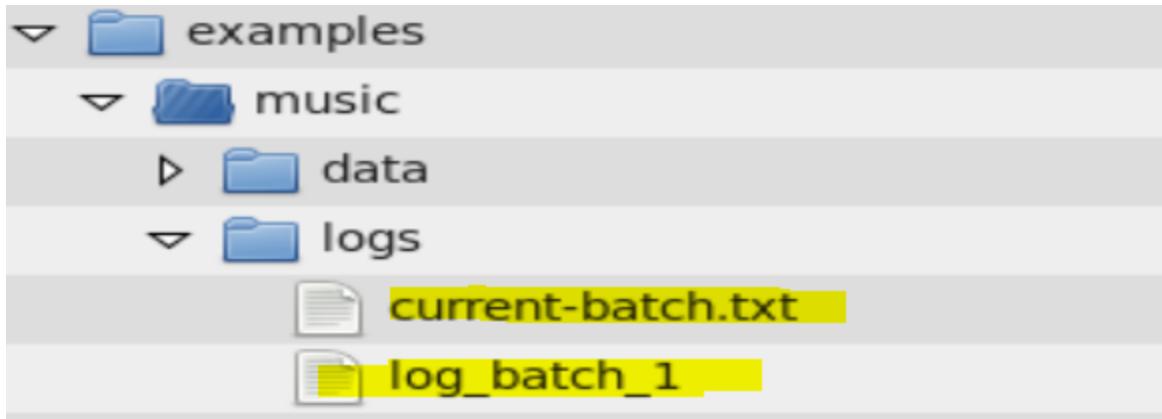
cat /home/acadgild/examples/music/logs/current-batch.txt
```

Output of the above script is validated through JPS Command

```
[acadgild@localhost music]$ ./music_project_master.sh
Preparing to execute python scripts to generate data...
Data Generated Successfully !
Starting the daemons...
After chmod
After batchid--> 1
This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh
18/08/06 11:48:17 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting namenodes on [localhost]
localhost: starting namenode, logging to /home/acadgild/install/hadoop/hadoop-2.6.5/logs/yarn-acadgild-namenode-localhost.localdomain.out
localhost: starting datanode, logging to /home/acadgild/install/hadoop/hadoop-2.6.5/logs/hadoop-acadgild-datanode-localhost.localdomain.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /home/acadgild/install/hadoop/hadoop-2.6.5/logs/hadoop-acadgild-secondarynamenode-localhost.localdomain.out
18/08/06 11:49:05 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
starting yarn daemons
starting resourcemanager, logging to /home/acadgild/install/hadoop/hadoop-2.6.5/logs/yarn-acadgild-resourcemanager-localhost.localdomain.out
localhost: starting nodemanager, logging to /home/acadgild/install/hadoop/hadoop-2.6.5/logs/yarn-acadgild-nodemanager-localhost.localdomain.out
localhost: starting zookeeper, logging to /home/acadgild/install/hbase/hbase-1.2.6/logs/hbase-acadgild-zookeeper-localhost.localdomain.out
starting master, logging to /home/acadgild/install/hbase/hbase-1.2.6/logs/hbase-acadgild-master-localhost.localdomain.out
starting regionserver, logging to /home/acadgild/install/hbase/hbase-1.2.6/logs/hbase-acadgild-1-regionserver-localhost.localdomain.out
starting historyserver, logging to /home/acadgild/install/hadoop/hadoop-2.6.5/logs/mapred-acadgild-historyserver-localhost.localdomain.out
17633 HMaster
6577 DataNode
7826 JobHistoryServer
6482 NameNode
7861 Jps
7733 HRegionServer
6726 SecondaryNameNode
7544 HQuorumPeer
6922 ResourceManager
7023 NodeManager
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 Mapping Complete !
Let us do data enrichment as per the requirement...
Data Enrichment Complete
Lets run some use cases now...
USE CASES COMPLETE !!
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost music]$
```

```
[acadgild@localhost music]$ jps
7633 HMaster
6577 DataNode
7826 JobHistoryServer
6482 NameNode
7733 HRegionServer
6726 SecondaryNameNode
7960 Jps
7544 HQuorumPeer
6922 ResourceManager
7023 NodeManager
[acadgild@localhost music]$
```

The **start-daemon.sh** script will check whether the current-batch.txt file is available in the logs folder or not. If not it will create the file and dump value '1' in that file and create LOGFILE with the current **batchid**.



2. This start-daemon.sh is getting used in music_project_master.sh , which is master file to execute the different scripts to data ingestion , activating daemons etc.

This is the master file which will be finally scheduled in Crontab for execution in every three hours.

```
# Create data
echo "Preparing to execute python scripts to generate data..."

rm -r /home/acadgild/examples/music/data/web
rm -r /home/acadgild/examples/music/data/mob
mkdir -p /home/acadgild/examples/music/data/web
mkdir -p /home/acadgild/examples/music/data/mob

python /home/acadgild/examples/music/generate_web_data.py
python /home/acadgild/examples/music/generate_mob_data.py

echo "Data Generated Successfully !"

# Call Stop start daemon scripts to start hadoop daemons

echo "Starting the daemons...."
#sh start-daemons.sh

# run jps commands to check the daemons

jps

echo "All hadoop daemons started !"

echo "Upload the look up tables now in Hbase...."

sh populate-lookup.sh

echo "Done with data population in look up tables !"
```

```

echo "Lets do some data formatting now...."
#sh dataformatting.sh
echo "data formatting complete !"
echo "Creating hive tables on top of hbase tables for data enrichment and filtering..."
#sh data_enrichment_filtering_schema.sh
echo "Hive table with Hbase Mapping Complete !"
echo "Let us do data enrichment as per the requirement..."
#sh data_enrichment.sh
echo "Data Enrichment Complete"

echo "Lets run some use cases now..."
#sh data_analysis.sh
echo "USE CASES COMPLETE !"

```

This script includes python scripts to create data.

3. Create data from mobile devices

```

from random import randint
from random import choice

file = open("/home/acadgild/examples/music/data/mob/file.txt", "w")
count = 20

while (count > 0):
    geo_cd_list=[ "A", "E", "AU", "AP", "U"]
    song_end_type_list=[ "0", "1", "2", "3"]
    timestamp_list=[ "1465230523", "1465130523", "1475130523", "1495130523"]
    start_ts_list=[ "1465230523", "1465130523", "1475130523", "1485130523"]
    end_ts_list=[ "1465230523", "1465130523", "1475130523", "1485130523"]

    if (count%15 == 0):
        user_id = ""
    else:
        user_id = "U" + str(randint(100,120))

    song_id = "S" + str(randint(200,210))

    if (count%11 == 0):
        artist_id = ""
    else:
        artist_id = "A" + str(randint(300,305))

    timestamp = choice(timestamp_list)
    start_ts = choice(start_ts_list)
    end_ts = choice(end_ts_list)

    if (count%12 == 0):
        geo_cd = ""
    else:
        geo_cd = choice(geo_cd_list)

```

```

station_id = "ST" + str(randint(400,415))
song_end_type = choice(song_end_type_list)
like = str(randint(0,1))
dislike = str(randint(0,1))

file.write("%s,%s,%s,%s,%s,%s,%s,%s\n" % (user_id, song_id, artist_id, timestamp, start_ts, end_ts,
geo_cd, station_id, song_end_type, like, dislike))

count = count-1

file.close()

```

4. Create data from Web

```

from random import randint
from random import choice

file = open("/home/acadgild/examples/music/data/web/file.xml", "w")
count = 20

file.write("<records>\n")

while (count > 0):
    geo_cd_list=[ "A", "E", "AU", "AP", "U"]
    song_end_type_list=[ "0","1","2","3"]
    timestamp_list=[ "2016-05-10 12:24:22", "2016-06-09 22:12:36", "2016-07-10 01:38:09", "2017-05-09 08:09:22"]
    start_ts_list=[ "2016-05-10 12:24:22", "2016-06-09 22:12:36", "2016-07-10 01:38:09", "2017-05-09 08:09:22"]
    end_ts_list=[ "2016-05-10 12:24:22", "2016-06-09 22:12:36", "2016-07-10 01:38:09", "2017-05-09 08:09:22"]

    if (count%15 == 0):
        user_id =""
    else:
        user_id = "U" + str(randint(100,120))

    song_id = "S" + str(randint(200,210))

    if (count%11 == 0):
        artist_id =""
    else:
        artist_id = "A" + str(randint(300,305))

    timestamp = choice(timestamp_list)
    start_ts = choice(start_ts_list)
    end_ts = choice(end_ts_list)

    if (count%12 == 0):
        geo_cd =""
    else:
        geo_cd = choice(geo_cd_list)

    timestamp = choice(timestamp_list)
    start_ts = choice(start_ts_list)
    end_ts = choice(end_ts_list)

    if (count%12 == 0):
        geo_cd =""
    else:
        geo_cd = choice(geo_cd_list)

    station_id = "ST" + str(randint(400,415))
    song_end_type = choice(song_end_type_list)
    like = str(randint(0,1))
    dislike = str(randint(0,1))

    file.write("<record>\n")
    file.write("<user_id>%s</user_id>\n" % (user_id))
    file.write("<song_id>%s</song_id>\n" % (song_id))
    file.write("<artist_id>%s</artist_id>\n" % (artist_id))
    file.write("<timestamp>%s</timestamp>\n" % (timestamp))
    file.write("<start_ts>%s</start_ts>\n" % (start_ts))
    file.write("<end_ts>%s</end_ts>\n" % (end_ts))
    file.write("<geo_cd>%s</geo_cd>\n" % (geo_cd))
    file.write("<station_id>%s</station_id>\n" % (station_id))
    file.write("<song_end_type>%s</song_end_type>\n" % (song_end_type))
    file.write("<like>%s</like>\n" % (like))
    file.write("<dislike>%s</dislike>\n" % (dislike))
    file.write("</record>\n")

    count = count-1

file.write("</records>")
file.close()

```

These scripts should be marked as executable and should display in green color if listed.
If it is not displayed as executable execute below command to change the permissions to execute

```
Chmod 774 /home/acadgild/examples/music/*
```

```
[acadgild@localhost music]$ chmod 775 /home/acadgild/examples/music/*
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost music]$ ls -l
total 48
-rwxrwxr-x. 1 acadgild acadgild 873 Aug  6 10:41 create_hive_hbase_lookup.hql
-rwxrwxr-x. 1 acadgild acadgild  84 Aug  6 10:41 create_hive_hbase_lookup.sh
drwxrwxr-x. 4 acadgild acadgild 4096 Aug  6 10:54 data
-rwxrwxr-x. 1 acadgild acadgild 299 Aug  6 10:41 data_enrichment_filtering_schema.sh
-rwxrwxr-x. 1 acadgild acadgild 1499 Aug  6 10:41 data_enrichment.sh
-rwxrwxr-x. 1 acadgild acadgild 402 Aug  6 10:41 dataformatting.sh
-rwxrwxr-x. 1 acadgild acadgild 1200 Aug  6 10:41 generate_mob_data.py
-rwxrwxr-x. 1 acadgild acadgild 1868 Aug  6 10:41 generate_web_data.py
-rwxrwxr-x. 1 acadgild acadgild 1256 Aug  6 10:50 music_project_master.sh
-rwxrwxr-x. 1 acadgild acadgild 1645 Aug  6 10:41 populate-lookup.sh
-rwxrwxr-x. 1 acadgild acadgild  773 Aug  6 10:30 start-daemons.sh
-rwxrwxr-x. 1 acadgild acadgild  337 Aug  6 10:41 user-artist.hql
[acadgild@localhost music]$
```

Output of Above Data Creation Script

```
[acadgild@localhost music]$ ./music_project_master.sh
Preparing to execute python scripts to generate data...
Data Generated Successfully !
Starting the daemons....
3627 Jps
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 Mapping Complete !
Let us do data enrichment as per the requirement...
Data Enrichment Complete
Lets run some use cases now...
USE CASES COMPLETE !!
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost music]$
```

data					
Name	Size	Type			
mob	1 item	folder			
file.txt	1.2 KB	plain text document	Mo		
web	1 item	folder			
file.xml	6.6 KB	XML document	Mo		

Section-4 –Data Ingestion, Formatting, Enrichment & Filtering, Analysis And Export

4.1 Stage – 1 – Data Ingestion

By using the “***populate-lookup.sh***” script we will create lookup tables in **Hbase**. These tables have to be used in,

Data formatting,

Data enrichment and

Analysis stage

Lookup Tables

Sl.no	Table Name	Description	Related File
1	station-geo-map	Contains mapping of a geo_cd with station_id	stn-geocd.txt
2	subscribed-users	Contains user_id , subscription_start_date and subscription_end_date . Contains details only for subscribed users	user-subscn.txt
3	song-artist-map	Contains mapping of song_id with artist_id Along with royalty associated with each play of the song	song-artist.txt
4	user-artist-map	Contains an array of artist_id(s) followed by a user_id	user-artist.txt

Table-1

“***populate-lookup.sh***” script

The “***populate-lookup.sh***” shell script creates the above 4 lookup tables in the Hbase and populate the data into the lookup tables from the dataset files.

In the below screen shots, we can see the create-lookup.sh scripts and the following screen shots shows the tables creation and population of the data in the Hbase. Also, the values loaded into the Hbase Tables are also shown, please see the below screen shots.

```

#!/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"

file="/home/acadgild/examples/music/lookupfiles/song-artist.txt"
while IFS= read -r line
do
  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"

  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"

  file="/home/acadgild/examples/music/lookupfiles/user-subscn.txt"
  while IFS= read -r line
  do
    userid=`echo $line | cut -d',' -f1`
    startdt=`echo $line | cut -d',' -f2`
    enddt=`echo $line | cut -d',' -f3`
    echo "put 'subscribed-users', '$userid', 'subscn:startdt', '$startdt'" | hbase shell
    echo "put 'subscribed-users', '$userid', 'subscn:enddt', '$enddt'" | hbase shell
  done <"$file"

#hive -f /home/acadgild/examples/music/user-artist.hql

```

Output of execution of this Script (This script is called through master script)

Create tables in Hbase

```
2018-08-06 12:05:10,188 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

create 'station-geo-map', 'geo'
0 row(s) in 3.6950 seconds

Hbase::Table - station-geo-map
2018-08-06 12:05:36,130 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

create 'subscribed-users', 'subscn'
0 row(s) in 2.1950 seconds

Hbase::Table - subscribed-users
2018-08-06 12:06:00,432 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

create 'song-artist-map', 'artist'
0 row(s) in 3.2430 seconds
```

```
hbase(main):001:0> list
TABLE
SparkHBasesTable
Transactions_Hbase
bulktable
clicks
song-artist-map
station-geo-map
subscribed-users
7 row(s) in 1.9550 seconds

=> ["SparkHBasesTable", "Transactions_Hbase", "bulktable", "clicks", "song-artist-map", "station-geo-map", "subscribed-users"]
hbase(main):002:0>
```

Populating the tables in Hbase

```
put 'station-geo-map', 'ST400', 'geo:geo_cd', 'A'
0 row(s) in 1.1880 seconds

2018-08-06 12:06:49,755 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

2018-08-06 12:27:37,367 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

put 'subscribed-users', 'U114', 'subscn:enddt', '1468130523'
0 row(s) in 0.0270 seconds

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 Mapping Complete !
Let us do data enrichment as per the requirement...
Data Enrichment Complete
Lets run some use cases now...
USE CASES COMPLETE !!
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost music]$
```

```
hbase(main):003:0> scan 'station-geo-map'
ROW
ST400
ST401
ST402
ST403
ST404
ST405
ST406
ST408
ST409
ST410
ST411
ST412
ST413
ST414
14 row(s) in 0.4250 seconds
          COLUMN+CELL
          column=geo:geo_cd, timestamp=1533537392680, value=A
          column=geo:geo_cd, timestamp=1533537415598, value=AU
          column=geo:geo_cd, timestamp=1533537439135, value=AP
          column=geo:geo_cd, timestamp=1533537462002, value=J
          column=geo:geo_cd, timestamp=1533537486062, value=E
          column=geo:geo_cd, timestamp=1533537523139, value=A
          column=geo:geo_cd, timestamp=1533537558397, value=E
          column=geo:geo_cd, timestamp=1533537589761, value=E
          column=geo:geo_cd, timestamp=1533537613642, value=E
          column=geo:geo_cd, timestamp=1533537637473, value=A
          column=geo:geo_cd, timestamp=1533537662163, value=A
          column=geo:geo_cd, timestamp=1533537686077, value=AP
          column=geo:geo_cd, timestamp=1533537709568, value=J
          column=geo:geo_cd, timestamp=1533537732167, value=E
```

```
hbase(main):004:0> scan 'song-artist-map'
ROW
S200
S201
S202
S203
S204
S205
S206
S207
S208
S209
10 row(s) in 0.2300 seconds
          COLUMN+CELL
          column=artist:artistid, timestamp=1533537754870, value=A300
          column=artist:artistid, timestamp=1533537778554, value=A301
          column=artist:artistid, timestamp=1533537802031, value=A302
          column=artist:artistid, timestamp=1533537825366, value=A303
          column=artist:artistid, timestamp=1533537848014, value=A304
          column=artist:artistid, timestamp=1533537871551, value=A301
          column=artist:artistid, timestamp=1533537895999, value=A302
          column=artist:artistid, timestamp=1533537919387, value=A303
          column=artist:artistid, timestamp=1533537942850, value=A304
          column=artist:artistid, timestamp=1533537966013, value=A305
```

```
hbase(main):005:0> scan 'subscribed-users'
ROW
U100
U100
U101
U101
U102
U102
U103
U103
U104
U104
U105
U105
U106
U106
U107
U107
U108
U108
U109
U109
U110
U110
U111
U111
U112
U112
U113
U113
U114
U114
15 row(s) in 0.3920 seconds
          COLUMN+CELL
          column=subscn:enddt, timestamp=1533538011867, value=1465130523
          column=subscn:startdt, timestamp=1533537989292, value=1465230523
          column=subscn:enddt, timestamp=1533538057056, value=1475130523
          column=subscn:startdt, timestamp=1533538034163, value=1465230523
          column=subscn:enddt, timestamp=1533538103501, value=1475130523
          column=subscn:startdt, timestamp=1533538080775, value=1465230523
          column=subscn:enddt, timestamp=1533538150964, value=1475130523
          column=subscn:startdt, timestamp=1533538127768, value=1465230523
          column=subscn:enddt, timestamp=1533538199062, value=1475130523
          column=subscn:startdt, timestamp=1533538174013, value=1465230523
          column=subscn:enddt, timestamp=1533538245359, value=1475130523
          column=subscn:startdt, timestamp=1533538222721, value=1465230523
          column=subscn:enddt, timestamp=1533538289879, value=1485130523
          column=subscn:startdt, timestamp=1533538267687, value=1465230523
          column=subscn:enddt, timestamp=1533538335427, value=1455130523
          column=subscn:startdt, timestamp=1533538312793, value=1465230523
          column=subscn:enddt, timestamp=1533538380942, value=1465230623
          column=subscn:startdt, timestamp=1533538358669, value=1465230523
          column=subscn:enddt, timestamp=1533538427546, value=1475130523
          column=subscn:startdt, timestamp=1533538404281, value=1465230523
          column=subscn:enddt, timestamp=1533538474669, value=1475130523
          column=subscn:startdt, timestamp=1533538451155, value=1465230523
          column=subscn:enddt, timestamp=1533538522002, value=1475130523
          column=subscn:startdt, timestamp=1533538498028, value=1465230523
          column=subscn:enddt, timestamp=1533538569399, value=1475130523
          column=subscn:startdt, timestamp=1533538545820, value=1465230523
          column=subscn:enddt, timestamp=1533538616647, value=1485130523
          column=subscn:startdt, timestamp=1533538593065, value=1465230523
          column=subscn:enddt, timestamp=1533538663315, value=1468130523
          column=subscn:startdt, timestamp=1533538639865, value=1465230523
```

The populate-lookup.sh also creates a lookup table “**users_artists**” in the HIVE, loading the data from the **user-artist.txt**, the below screen shot shows that the table has been created in the HIVE.

```
Logging initialized using configuration in jar:file:/home/acadgild/install/hive/apache-hive-2.3.2-bin/lib/hive-common-2.3.2.jar!/hive-log4j2.properties Async: true
OK
Time taken: 20.453 seconds
OK
Time taken: 0.098 seconds
OK
Time taken: 2.846 seconds
Loading data to table project.users_artists
OK
Time taken: 5.798 seconds
You have new mail in /var/spool/mail/acadgild
```

```
hive> show databases;
OK
acadgilddb
custom
default
project
transactions
Time taken: 0.164 seconds, Fetched: 5 row(s)
hive> use project
      > ;
OK
Time taken: 0.112 seconds
hive> show tables;
OK
Time taken: 0.177 seconds
hive> show tables;
OK
users_artists
Time taken: 0.38 seconds, Fetched: 1 row(s)
hive> █
```

```
hive> select * from users_artists;
OK
U100  ["A300","A301","A302"]
U101  ["A301","A302"]
U102  ["A302"]
U103  ["A303","A301","A302"]
U104  ["A304","A301"]
U105  ["A305","A301","A302"]
U106  ["A301","A302"]
U107  ["A302"]
U108  ["A300","A303","A304"]
U109  ["A301","A303"]
U110  ["A302","A301"]
U111  ["A303","A301"]
U112  ["A304","A301"]
U113  ["A305","A302"]
U114  ["A300","A301","A302"]
Time taken: 7.431 seconds, Fetched: 15 row(s)
hive> █
```

Now we need to link these lookup tables in hive using the Hbase Storage Handler.
With the help of “**data_enrichment_filtering_schema.sh**” file we will create hive tables on the top of Hbase tables using “**create_hive_hbase_lookup.hql**”.

Creating Hive Tables on the top of Hbase:

In this section 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.

Run the script: **./data_enrichment_filtering_schema.sh, (through master script only)**

The script will run the “**create_hive_hbase_lookup.hql**” which will create the HIVE external tables with the help of **Hbase storage handler & SerDe properties**. The hive external tables will match the columns of **Hbase** tables to **HIVE** tables.

```
#!/bin/bash
batchid=`cat /home/acadgild/project/logs/current-batch.txt`
LOGFILE=/home/acadgild/project/logs/log_batch_$batchid
echo "Creating hive tables on top of hbase tables for data enrichment and filtering..." >> $LOGFILE
hive -f /home/acadgild/project/scripts/create_hive_hbase_lookup.hql
```

create_hive_hbase_lookup.hql

```
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'
with serdeproperties
("hbase.columns.mapping"=:key,artist:artistid")
tblproperties("hbase.table.name"="song-artist-map");
```

```
[acadgild@localhost music]$ ./music_project_master.sh
Preparing to execute python scripts to generate data...
7633 HMaster
6577 DataNode
7826 JobHistoryServer
6482 NameNode
7733 HRegionServer
6726 SecondaryNameNode
7544 HQuorumPeer
14794 RunJar
6922 ResourceManager
16366 Jps
7023 NodeManager
All hadoop daemons started !
Creating hive tables on top of hbase tables for data enrichment and filtering...
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/acadgild/install/hive/apache-hive-2.3.2-bin/lib/log4j-slf4j-impl-2.6.2.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.apache.logging.slf4j.Log4jLoggerFactory]

Logging initialized using configuration in jar:file:/home/acadgild/install/hive/apache-hive-2.3.2-bin/lib/hive-common-2.3.2.jar!/hive-log4j2.properties A
sync: true
OK
Time taken: 19.684 seconds
OK
Time taken: 9.729 seconds
OK
Time taken: 0.988 seconds
OK
Time taken: 1.023 seconds
Hive table with Hbase Mapping Complete !
USE CASES COMPLETE !!
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost music]$
```

```
hive> show tables;
OK
song_artist_map
station_geo_map
subscribed_users
users_artists
Time taken: 0.136 seconds, Fetched: 4 row(s)
hive>
```

```
hive> select * from station_geo_map;
OK
ST400    A
ST401    AU
ST402    AP
ST403    J
ST404    E
ST405    A
ST406    AU
ST408    E
ST409    E
ST410    A
ST411    A
ST412    AP
ST413    J
ST414    E
Time taken: 1.497 seconds, Fetched: 14 row(s)
hive>
```

```
hive> select * from song_artist_map;
OK
S200      A300
S201      A301
S202      A302
S203      A303
S204      A304
S205      A301
S206      A302
S207      A303
S208      A304
S209      A305
Time taken: 0.953 seconds, Fetched: 10 row(s)
hive> █
```

```
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
U114      1465230523      1468130523
Time taken: 1.08 seconds, Fetched: 15 row(s)
hive> █
```

4.2 Stage – 2 - Data Formatting

In this stage we are merging the data coming from both **web** applications and **mobile** applications and create a common table for analyzing purpose and create partitioned data based on **batched**.

Run the script: *./dataformatting.sh (Through master script only)*

```
#!/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 formatting..." >> $LOGFILE
spark-submit --packages com.databricks:spark-xml_2.10:0.4.1 \
--class DataFormatting \
--master local[2] \
/home/acadgild/examples/music/MusicDataAnalysis/target/scala-2.11/musicdataanalysis_2.11-1.0.jar $batchid
```

Here we can see DataFormatting class will be submitted by Spark.

But prior to execution of above – get all the dependent jars downloaded for DataFormatting.scala using **sbt -v package** command (might take 30mins). This command downloads all the JARs required for the dependencies mentioned in the **buid.sbt** file.

Execution of sbt -v package

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

Getting org.scala-sbt sbt 1.0.4 (this may take some time)...
```

```

Getting org.scala-sbt sbt 1.0.4 (this may take some time)...
downloading https://repo1.maven.org/maven2/org/scala-sbt/util-position_2.12/1.0.3/util-position_2.12-1.0.3.jar ...
[SUCCESSFUL ] org.scala-sbt#util-position_2.12;1.0.3!util-position_2.12.jar (4252ms)
downloading https://repo1.maven.org/maven2/com/eed3si9n/json-new-core_2.12/0.8.0/json-new-core_2.12-0.8.0.jar ...
[SUCCESSFUL ] com.eed3si9n#json-new-core_2.12;0.8.0!json-new-core_2.12.jar (1528ms)
downloading https://repo1.maven.org/maven2/com/eed3si9n/shaded-scalajson_2.12/1.0.0-M4/shaded-scalajson_2.12-1.0.0-M4.jar ...
[SUCCESSFUL ] com.eed3si9n#shaded-scalajson_2.12;1.0.0-M4!shaded-scalajson_2.12.jar (745ms)
downloading https://repo1.maven.org/maven2/org/spire-math/jawn-parser_2.12/0.10.4/jawn-parser_2.12-0.10.4.jar ...
[SUCCESSFUL ] org.spire-math#jawn-parser_2.12;0.10.4!jawn-parser_2.12.jar (701ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/completion_2.12;1.0.4!completion_2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#completion_2.12;1.0.4!completion_2.12.jar (788ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/task-system_2.12/1.0.4/task-system_2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#task-system_2.12;1.0.4!task-system_2.12.jar (773ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/tasks_2.12/1.0.4/tasks_2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#tasks_2.12;1.0.4!tasks_2.12.jar (748ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/testing_2.12/1.0.4/testing_2.12-1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#testing_2.12;1.0.4!testing_2.12.jar (747ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/util-tracking_2.12/1.0.3/util-tracking_2.12-1.0.3.jar ...
[SUCCESSFUL ] org.scala-sbt#util-tracking_2.12;1.0.3!util-tracking_2.12.jar (725ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/zinc-classpath_2.12/1.0.5/zinc-classpath_2.12-1.0.5.jar ...
[SUCCESSFUL ] org.scala-sbt#zinc-classpath_2.12;1.0.5!zinc-classpath_2.12.jar (686ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/zinc-apiinfo_2.12/1.0.5/zinc-apiinfo_2.12-1.0.5.jar ...
[SUCCESSFUL ] org.scala-sbt#zinc-apiinfo_2.12;1.0.5!zinc-apiinfo_2.12.jar (791ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/zinc-ivy-integration_2.12/1.0.5/zinc-ivy-integration_2.12-1.0.5.jar ...
[SUCCESSFUL ] org.scala-sbt#zinc-ivy-integration_2.12;1.0.5!zinc-ivy-integration_2.12.jar (699ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/zinc_2.12/1.0.5/zinc_2.12-1.0.5.jar ...
[SUCCESSFUL ] org.scala-sbt#zinc_2.12;1.0.5!zinc_2.12.jar (727ms)
downloading https://repo1.maven.org/maven2/jline/jline/2.14.4/jline-2.14.4.jar ...
[SUCCESSFUL ] jline#jline;2.14.4!jline.jar (786ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/util-control_2.12/1.0.3/util-control_2.12-1.0.3.jar ...
[SUCCESSFUL ] org.scala-sbt#util-control_2.12;1.0.3!util-control_2.12.jar (708ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/util-interface_1.0.3/util-interface-1.0.3.jar ...
[SUCCESSFUL ] org.scala-sbt#util-interface_1.0.3!util-interface.jar (708ms)
downloading https://repo1.maven.org/maven2/com/lmax/disruptor/3.3.6/disruptor-3.3.6.jar ...
[SUCCESSFUL ] com.lmax#disruptor;3.3.6!disruptor.jar (711ms)
downloading https://repo1.maven.org/maven2/org/scala-lang/scala-reflect_2.12.4/scala-reflect_2.12.4.jar ...
[SUCCESSFUL ] org.scala-lang#scala-reflect_2.12.4!scala-reflect.jar (3457ms)
downloading https://repo1.maven.org/maven2/org/scala-lang/scala-compiler_2.12.4/scala-compiler_2.12.4.jar ...
[SUCCESSFUL ] org.scala-lang#scala-compiler_2.12.4!scala-compiler.jar (9924ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/util-cache_2.12/1.0.3/util-cache_2.12-1.0.3.jar ...
[SUCCESSFUL ] org.scala-sbt#util-cache_2.12;1.0.3!util-cache_2.12.jar (846ms)
downloading https://repo1.maven.org/maven2/com/eed3si9n/json-new-murmurhash_2.12/0.8.0/json-new-murmurhash_2.12-0.8.0.jar ...
[SUCCESSFUL ] com.eed3si9n#json-new-murmurhash_2.12;0.8.0!json-new-murmurhash_2.12.jar (679ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/test-agent_1.0.4/test-agent_1.0.4.jar ...
[SUCCESSFUL ] org.scala-sbt#test-agent_1.0.4!test-agent.jar (697ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/test-interface_1.0/test-interface_1.0.jar ...
[SUCCESSFUL ] org.scala-sbt#test-interface_1.0!test-interface.jar (702ms)
downloading https://repo1.maven.org/maven2/org/scala-sbt/zinc-classfile_2.12/1.0.5/zinc-classfile_2.12-1.0.5.jar ...
[SUCCESSFUL ] org.scala-sbt#zinc-classfile_2.12;1.0.5!zinc-classfile_2.12.jar (731ms)
downloading https://repo1.maven.org/maven2/com.jcraft/jsch/0.1.46/jsch-0.1.46.jar ...
[SUCCESSFUL ] com.jcraft#jsch;0.1.46!jsch.jar (884ms)
downloading https://repo1.maven.org/maven2/com/eed3si9n/gigahorse-okhttp_2.12/0.3.0/gigahorse-okhttp_2.12-0.3.0.jar ...
[SUCCESSFUL ] com.eed3si9n#gigahorse-okhttp_2.12;0.3.0!gigahorse-okhttp_2.12.jar (714ms)
downloading https://repo1.maven.org/maven2/com/squareup.okhttp3 okhttp3/okhttp-urlconnection_3.7.0/okhttp-urlconnection_3.7.0.jar ...
[SUCCESSFUL ] com.squareup.okhttp3#okhttp-urlconnection_3.7.0!okhttp-urlconnection_3.7.0.jar (701ms)
downloading https://repo1.maven.org/maven2/com/eed3si9n/gigahorse-core_2.12/0.3.0/gigahorse-core_2.12-0.3.0.jar ...
[SUCCESSFUL ] com.eed3si9n#gigahorse-core_2.12;0.3.0!gigahorse-core_2.12.jar (841ms)
downloading https://repo1.maven.org/maven2/com/squareup.okhttp3 okhttp3/3.7.0/okhttp-3.7.0.jar ...

[info]  !SUCCESSFUL | com.jolbox#bonecp;0.8.0.RELEASE!bonecp.jar(bundle) (652ms)
[info] downloading https://repo1.maven.org/maven2/org/apache/derby/derby/derby/10.10.2.0/derby-10.10.2.0.jar ...
[info]  !SUCCESSFUL | org.apache.derby#derby/10.10.2.0!derby.jar (2229ms)
[info] downloading https://repo1.maven.org/maven2/org/datanucleus/datanucleus-api-jdo/3.2.6/datanucleus-api-jdo-3.2.6.jar ...
[info]  !SUCCESSFUL | org.datanucleus#datanucleus-api-jdo;3.2.6!datanucleus-api-jdo.jar (760ms)
[info] downloading https://repo1.maven.org/maven2/org/datanucleus/datanucleus-rdbms/3.2.9/datanucleus-rdbms-3.2.9.jar ...
[info]  !SUCCESSFUL | org.datanucleus#datanucleus-rdbms;3.2.9!datanucleus-rdbms.jar (1522ms)
[info] downloading https://repo1.maven.org/maven2/commons-pool/commons-pool_1.5.4/commons-pool-1.5.4.jar ...
[info]  !SUCCESSFUL | commons-pool#commons-pool_1.5.4!commons-pool.jar (639ms)
[info] downloading https://repo1.maven.org/maven2/commons-dbcp/commons-dbcp_1.4/commons-dbcp-1.4.jar ...
[info]  !SUCCESSFUL | commons-dbcp#commons-dbcp;1.4!commons-dbcp.jar (650ms)
[info] downloading https://repo1.maven.org/maven2/javax/do/dojo/api/3.0.1/do/api-3.0.1.jar ...
[info]  !SUCCESSFUL | javax.do#do/api/3.0.1!do/api.jar (677ms)
[info] downloading https://repo1.maven.org/maven2/javax/transaction/jta/1.1/jta-1.1.jar ...
[info]  !SUCCESSFUL | javax.transaction#jta;1.1!jta.jar (674ms)
[info] downloading https://repo1.maven.org/maven2/org/apache/calcite/calcite-linq4j/1.2.0-incubating/calcite-linq4j-1.2.0-incubating.jar ...
[info]  !SUCCESSFUL | org.apache.calcite#calcite-linq4j;1.2.0-incubating!calcite-linq4j.jar (833ms)
[info] downloading https://repo1.maven.org/maven2/net/hydromatic/eigenbase-properties/1.1.5/eigenbase-properties-1.1.5.jar ...
[info]  !SUCCESSFUL | net.hydromatic#eigenbase-properties;1.1.5!eigenbase-properties.jar(bundle) (600ms)
[info] downloading https://repo1.maven.org/maven2/org/apache/httpcomponents/httpcore/4.4.4/httpcore-4.4.4.jar ...
[info]  !SUCCESSFUL | org.apache.httpcomponents#httpcore;4.4.4!httpcore.jar (745ms)
[info] downloading https://repo1.maven.org/maven2/commons-logging/commons-logging/1.2/commons-logging-1.2.jar ...
[info]  !SUCCESSFUL | commons-logging#commons-logging;1.2!commons-logging.jar (621ms)
[info] downloading https://repo1.maven.org/maven2/jline/jline/2.12.1/jline-2.12.1.jar ...
[info]  !SUCCESSFUL | jline#jline;2.12.1!jline.jar (669ms)
[info] Done updating.
[info] Compiling 1 Scala source to /home/acadgild/examples/music/MusicDataAnalysis/target/scala-2.11/classes ...
[info] Non-compiled module 'compiler-bridge_2.11' for Scala 2.11.8. Compiling...
[info] Compilation completed in 53.451s.
[warn] there was one deprecation warning; re-run with -deprecation for details
[warn] one warning found
[info] Done compiling.
[info] Packaging /home/acadgild/examples/music/MusicDataAnalysis/target/scala-2.11/musicdataanalysis_2.11-1.0.jar ...
[info] Done packaging.
[success] Total time: 562 s, completed Aug 7, 2018 11:34:37 AM
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost MusicDataAnalysis]$ 
```

Dataformatting.scala

```

import org.apache.spark.{SparkConf, SparkContext}
import org.apache.spark.sql

object DataFormatting {
  def main(args: Array[String]): Unit = {
    val conf = new SparkConf().setAppName("Data Formatting")
    val sc = new SparkContext(conf)
    val sqlContext = new org.apache.spark.sql.hive.HiveContext(sc)
    val batchId = args(0)
    val create_hive_table = """CREATE TABLE IF NOT EXISTS project.formatted_input
      (
        User_id STRING,
        Song_id STRING,
        Artist_id STRING,
        Timestamp STRING,
        Start_ts STRING,
        End_ts STRING,
        Geo_cd STRING,
        Station_id STRING,
        Song_end_type INT,
        Like INT,
        Dislike INT
      )
      PARTITIONED BY
      (batchid INT)
      ROW FORMAT DELIMITED
      FIELDS TERMINATED BY ','
    """
    val load_mob_data = s"""LOAD DATA LOCAL INPATH '/home/acadgild/project/data/mob/file.txt'
      INTO TABLE project.formatted_input PARTITION (batchid='$batchId')"""
    val load_web_data = s"""INSERT INTO project.formatted_input
      PARTITION(batchid='$batchId')
      SELECT user_id,
      song_id,
      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,
      unix_timestamp(end_ts,'yyyy-MM-dd HH:mm:ss') AS end_ts,
      geo_cd,
      station_id,
      song_end_type,
      like,
      dislike
      FROM web_data
    """
    try {
      val xmlData = sqlContext.read.format(source = "com.databricks.spark.xml").option("rowTag", "record").load(path = "/home/acadgild/project/data/web/file.xml")
      xmlData.createOrReplaceTempView(viewName = "web_data")

      sqlContext.sql(create_hive_table)
      sqlContext.sql(load_mob_data)
      sqlContext.sql(load_web_data)
    }
    catch{
      case e: Exception=>e.printStackTrace()
    }
  }
}

```

Make sure that Hive metastore service (use command – “**hive –service metastore**”) is running and then we are executing ***dataformatting.sh***

We can submit your Spark application to a Spark deployment environment for execution, kill or request status of Spark applications.

Dataformatting.sh

```
#!/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 formatting..." >> $LOGFILE

spark-submit --packages com.databricks:spark-xml_2.10:0.4.1 \
--class DataFormatting \
--master local[2] \
/home/acadgild/examples/music/MusicDataAnalysis/target/scala-2.11/musicdataanalysis_2.11-1.0.jar $batchid
```

Execution of hive service metastore

Make sure mysql is also started

```
[acadgild@localhost ~]$ sudo service mysqld start [ OK ]
Starting mysqld:
[acadgild@localhost ~]$ mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 8
Server version: 8.0.3-rc-log MySQL Community Server (GPL)

Copyright (c) 2000, 2017, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement
mysql> █
```

The below screenshots represent the flow in the terminal.

```
[acadgild@localhost music]$ ./dataformatting.sh
Ivy Default Cache set to: /home/acadgild/.ivy2/cache
The jars for the packages stored in: /home/acadgild/.ivy2/jars
:: loading settings :: url = jar:file:/home/acadgild/install/spark/spark-2.2.1-bin-hadoop2.7/jars/ivy-2.4.0.jar!/org/apache/ivy/core/settings/ivysettings.xml
com.databricks#spark-xml_2.10 added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent;1.0
  confs: [default]
    found com.databricks#spark-xml_2.10;0.4.1 in central
:: resolution report :: resolve 1250ms :: artifacts dl 23ms
  :: modules in use:
    com.databricks#spark-xml_2.10;0.4.1 from central in [default]
-----
|   conf      | number| modules          ||  artifacts
|             |       | search|dwnlded|evicted|| number|dwnlded|
| default     |   1   |   0   |   0   |   0   ||   1   |   0   |
-----
:: retrieving :: org.apache.spark#spark-submit-parent
  confs: [default]
    0 artifacts copied, 1 already retrieved (0kB/47ms)
18/08/07 12:00:13 INFO spark.SparkContext: Running Spark version 2.2.1
18/08/07 12:00:14 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
18/08/07 12:00:14 WARN util.Utils: Your hostname, localhost.localdomain resolves to a loopback address: 127.0.0.1; using 192.168.0.4 instead (on interface eth12)
18/08/07 12:00:14 WARN util.Utils: Set SPARK LOCAL IP if you need to bind to another address
18/07/24 21:11:42 INFO execution.SparkSqlParser: Parsing command: CREATE TABLE IF NOT EXISTS project.formatted_input
(
  User_id STRING,
  Song_id STRING,
  Artist_id STRING,
  Timestamp STRING,
  Start_ts STRING,
  End_ts STRING,
  Geo_cd STRING,
  Station_id STRING,
  Song_end_type INT,
  Like INT,
  Dislike INT
)
PARTITIONED BY
(batchid INT)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
18/07/24 21:11:43 INFO parser.CatalystSqlParser: Parsing command: arrays<string>
18/07/24 21:11:45 INFO execution.SparkSqlParser: Parsing command: LOAD DATA LOCAL INPATH 'file:///home/acadgild/examples/music/data/mob/file.txt'
INTO TABLE project.formatted_input PARTITION (batchid='1')
18/07/24 21:11:46 INFO parser.CatalystSqlParser: Parsing command: int
```

New table is created in HIVE

```
hive> show tables;
OK
formatted_input
song_artist_map
station_geo_map
subscribed_users
users_artists
Time taken: 0.08 seconds, Fetched: 5 row(s)
hive>
```

Records from the table

hive> select * from formatted_input;												
U103	S201	A302	1465230523	1485130523	1485130523	AP	ST411	0	0	1	1	
U102	S200	A300	1475130523	1465130523	1485130523	AP	ST412	3	0	1	1	
U112	S202	A305	1495130523	1465230523	1465130523	AU	ST404	1	1	0	1	
U117	S204	A300	1475130523	1475130523	1485130523	U	ST407	1	0	0	1	
U108	S208	A306	1465230523	1475130523	1465130523	U	ST401	2	0	1	1	
	S202	A305	1475130523	1465230523	1465130523	AP	ST410	0	0	0	1	
U116	S203	A305	1465130523	1465230523	1485130523	E	ST414	1	0	0	1	
U118	S210	A301	1465130523	1485130523	1475130523	E	ST413	2	1	0	1	
U112	S205	A301	1465230523	1485130523	1475130523		ST401	2	1	1	1	
U109	S207		1465230523	1485130523	1465130523	E	ST413	3	1	1	1	
U112	S201	A303	1495130523	1485130523	1465230523	AP	ST411	1	0	1	1	
U113	S206	A302	1495130523	1485130523	1465130523	E	ST400	0	1	0	1	
U112	S200	A300	1475130523	1465230523	1465130523	A	ST414	2	0	0	1	
U114	S210	A300	1465130523	1475130523	1475130523	E	ST409	3	0	0	1	
U109	S207	A303	1465130523	1485130523	1475130523	E	ST402	1	0	0	1	
U115	S209	A301	1495130523	1485130523	1475130523	AP	ST406	2	0	1	1	
U112	S203	A302	1465230523	1465130523	1465130523	AP	ST413	1	1	1	1	
U116	S206	A300	1465230523	1485130523	1465230523	AU	ST410	3	0	0	1	
U107	S210	A302	1465130523	1465230523	1465130523	A	ST405	3	0	0	1	
U103	S206	A305	1465130523	1465130523	1475130523	AU	ST407	3	1	0	1	
U114	S206	A304	1462863262	1468094889	1462863262	AP	ST403	1	1	0	1	
U104	S205	A305	1462863262	1462863262	1468094889	AU	ST405	2	1	1	1	
U206	S206	A302	1462863262	1465490556	1494297562	AU	ST400	3	0	0	1	
U111	S201	A305	1494297562	1465490556	1494297562	AP	ST405	3	0	0	1	
U118	S208	A302	1468094889	1462863262	1468094889	U	ST407	0	0	1	1	

In the above screenshot we can see the formatted input data with some null values in **user_id**, **artist_id** and **geo_cd** columns which we will fill the enrichment script based on rules of enrichment for **artist_id** and **geo_cd** only. We will get neglect **user_id** because they didn't mention anything about **user_id** for enrichment purpose.

4.3 Data Enrichment and filtering

In this stage, 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.

Rules for 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**.

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.

data_enrichment.sh

```

#!/bin/bash

batchid=`cat /home/acadgild/project/logs/current-batch.txt`
LOGFILE=/home/acadgild/examples/music/logs/log_batch_$batchid
VALIDDIR=/home/acadgild/examples/music/processed_dir/valid/batch_$batchid
INVALIDDIR=/home/acadgild/examples/music/processed_dir/invalid/batch_$batchid

echo "Running script for data enrichment and filtering..." >> $LOGFILE

spark-submit --class DataEnrichment \
--master local[2] \
--jars /usr/local/hive/lib/hive-hbase-handler-2.1.1.jar,/usr/local/hive/lib/hbase-client-1.1.1.jar,/usr/local/hive/lib/hbase-common-1.1.1.jar,/usr/local/hive/lib/hbase-hadoop-compat-1.1.1.jar,/usr/local/hive/lib/hbase-server-1.1.1.jar,/usr/local/hive/lib/hbase-protocol-1.1.1.jar,/usr/local/hive/lib/zookeeper-3.4.6.jar,/usr/local/hive/lib/guava-14.0.1.jar,/usr/local/hive/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

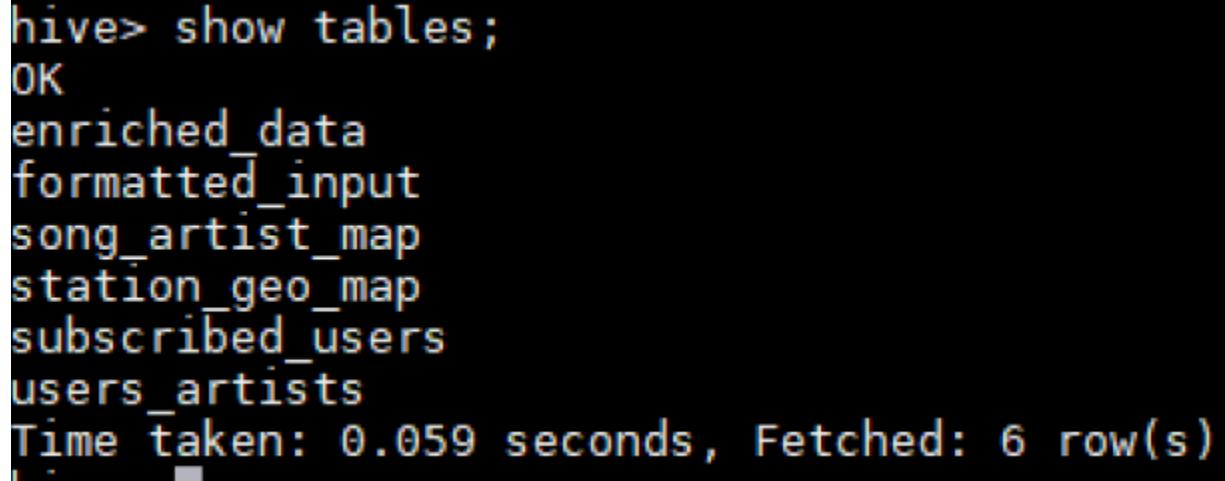
hadoop fs -get /user/hive/warehouse/project.db/enriched_data/batchid=$batchid/status=pass/* $VALIDDIR
hadoop fs -get /user/hive/warehouse/project.db/enriched_data/batchid=$batchid/status=fail/* $INVALIDDIR

echo "Deleting older valid and invalid records from local file system..." >> $LOGFILE

find /home/acadgild/examples/music/processed_dir/ -mtime +7 -exec rm {} \;

```

Once running the above script in terminal we would get a new table created in Hive. Please find the below screenshots for the Hive screens.



```

hive> show tables;
OK
enriched_data
formatted_input
song_artist_map
station_geo_map
subscribed_users
users_artists
Time taken: 0.059 seconds, Fetched: 6 row(s)

```

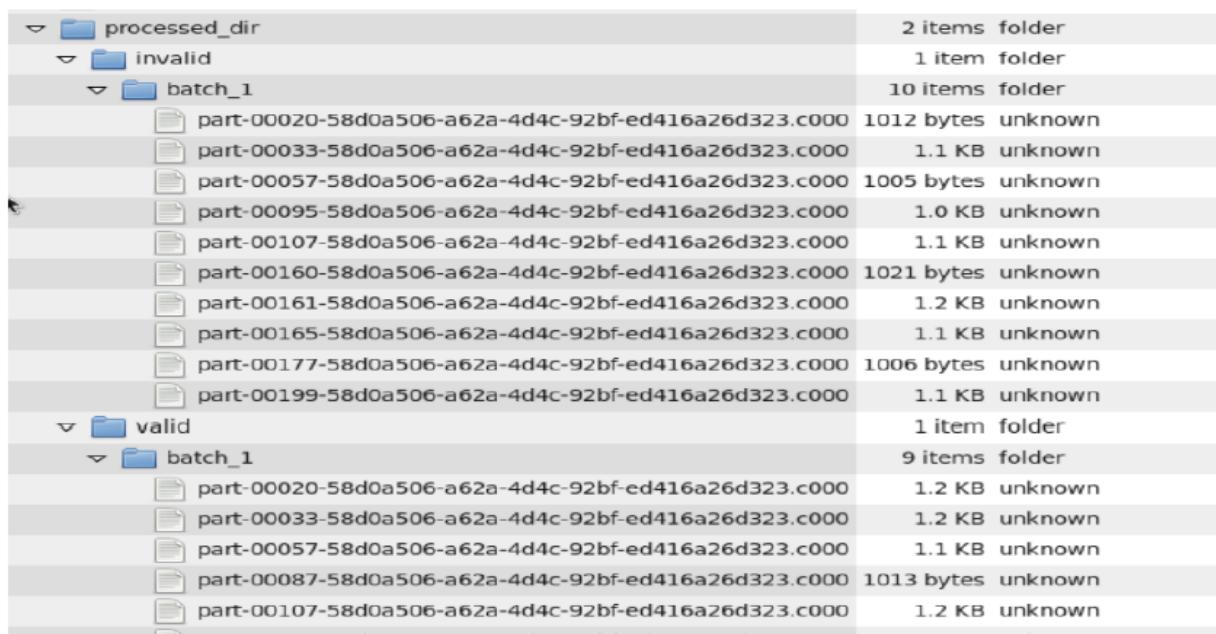
```

hive> select * from enriched_data;
OK
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| U113 | S201 | A301 | 1462863262 | 1465490556 | 1465490556 | NULL | ST415 | 1 | 0 | 1 | 1 | 1 | fail |
| NULL | S207 | A303 | 1462863262 | 1494297562 | 1468094889 | NULL | ST414 | 0 | 0 | 0 | 1 | 1 | fail |
| U109 | S207 | A303 | 1465230523 | 1485130523 | 1465130523 | J | ST413 | 3 | 1 | 1 | 1 | 1 | fail |
| S202 | A302 | 1475130523 | 1465230523 | 1465130523 | A | ST410 | 0 | 0 | 0 | 1 | 1 | fail |
| U115 | S209 | NULL | 1495130523 | 1485130523 | 1475130523 | AU | ST406 | 2 | 0 | 1 | 1 | 1 | fail |
| U108 | S206 | A302 | 1465490556 | 1462863262 | 1462863262 | J | ST403 | 1 | 1 | 1 | 1 | 1 | fail |
| U114 | S206 | A302 | 1468094889 | 1494297562 | 1494297562 | E | ST409 | 3 | 0 | 0 | 1 | 1 | fail |
| U101 | S208 | A304 | 1494297562 | 1468094889 | 1494297562 | NULL | ST415 | 1 | 1 | 0 | 1 | 1 | fail |
| U107 | S210 | NULL | 1465130523 | 1465230523 | 1465130523 | A | ST405 | 3 | 0 | 0 | 1 | 1 | fail |
| U114 | S210 | NULL | 1465130523 | 1475130523 | 1475130523 | E | ST409 | 3 | 0 | 0 | 1 | 1 | fail |
| U119 | S210 | NULL | 1494297562 | 1465490556 | 1465490556 | A | ST410 | 1 | 0 | 1 | 1 | 1 | fail |
| U118 | S210 | NULL | 1465130523 | 1485130523 | 1475130523 | J | ST413 | 2 | 1 | 0 | 1 | 1 | fail |
| U108 | S210 | NULL | 1468094889 | 1465490556 | 1462863262 | J | ST413 | 3 | 0 | 0 | 1 | 1 | fail |
| U104 | S205 | A301 | 1462863262 | 1462863262 | 1468094889 | A | ST405 | 2 | 1 | 1 | 1 | 1 | fail |
| U112 | S205 | A301 | 1465230523 | 1485130523 | 1475130523 | AU | ST401 | 2 | 1 | 1 | 1 | 1 | fail |
| U112 | S200 | A300 | 1475130523 | 1465230523 | 1465230523 | NULL | ST414 | 2 | 0 | 0 | 1 | 1 | fail |
| U116 | S203 | A303 | 1465130523 | 1485130523 | 1485130523 | NULL | ST414 | 1 | 0 | 0 | 1 | 1 | fail |
| U112 | S203 | A303 | 1465230523 | 1465130523 | 1465130523 | J | ST413 | 1 | 1 | 1 | 1 | 1 | fail |
| U111 | S201 | A301 | 1494297562 | 1465490556 | 1494297562 | A | ST405 | 3 | 0 | 0 | 1 | 1 | pass |
| U103 | S201 | A301 | 1465230523 | 1485130523 | 1485130523 | A | ST411 | 0 | 0 | 1 | 1 | 1 | pass |
| U112 | S201 | A301 | 1495130523 | 1485130523 | 1465230523 | A | ST411 | 1 | 0 | 1 | 1 | 1 | pass |
| U120 | S201 | A301 | 1494297562 | 1462863262 | 1465490556 | J | ST413 | 1 | 1 | 0 | 1 | 1 | pass |
| U109 | S207 | A303 | 1465130523 | 1485130523 | 1475130523 | AP | ST402 | 1 | 0 | 0 | 1 | 1 | pass |
| U119 | S207 | A303 | 1468094889 | 1494297562 | 1468094889 | J | ST403 | 1 | 0 | 0 | 1 | 1 | pass |
| U100 | S207 | A303 | 1462863262 | 1462863262 | 1465490556 | A | ST410 | 0 | 1 | 0 | 1 | 1 | pass |
| U112 | S202 | A302 | 1495130523 | 1465230523 | 1465130523 | E | ST404 | 1 | 1 | 0 | 1 | 1 | pass |
| U117 | S202 | A302 | 1465490556 | 1465490556 | 1465490556 | E | ST409 | 0 | 0 | 0 | 1 | 1 | pass |
| U117 | S204 | A304 | 1475130523 | 1475130523 | 1465130523 | AP | ST407 | 1 | 0 | 0 | 1 | 1 | pass |
| U113 | S206 | A302 | 1495130523 | 1485130523 | 1465130523 | A | ST400 | 0 | 1 | 0 | 1 | 1 | pass |
| U104 | S206 | A302 | 1462863262 | 1465490556 | 1494297562 | A | ST400 | 3 | 0 | 0 | 1 | 1 | pass |
| U114 | S206 | A302 | 1462863262 | 1468094889 | 1462863262 | J | ST403 | 1 | 1 | 0 | 1 | 1 | pass |

```

In the above screenshot we can see the records are categorised as “pass” or “fail”.

All the “**pass**” records are accumulated in “**processed_dir/valid/batch_1**” and “**fail**” records are accumulated in “**processed_dir/invalid/batch_1**” (refer below screenshot)



4.4 Data Analysis

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

All the Spark SQL are captured in **DataAnalysis.scala** which will be internally triggered by **data_analysis.sh**.

```
#!/bin/bash
batchid=`cat /home/acadgild/example/music/logs/current-batch.txt`
LOGFILE=/home/acadgild/example/music/logs/log_batch_$batchid

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

spark-submit --class DataAnalysis --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-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/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 -master local[*] /home/acadgild/example/music/MusicDataAnalysis/target/scala-2.11/musicedataanalysis_2.11-1.0.jar $batchid
```

DataAnalysis.scala

```
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)

    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 = """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 likes1
GROUP BY station_id
ORDER BY total_distinct_songs_played DESC
LIMIT 10"""

    val create_users_behaviour = """CREATE TABLE IF NOT EXISTS users_behaviour
(
user_type STRING,
duration INT
)
PARTITIONED BY (batchid INT)
ROW FORMAT DELIMITED
```

```

STORED AS TEXTFILE"""
|   val load_usersBehaviour = """INSERT OVERWRITE TABLE usersBehaviour
PARTITION(batchid='$batchId')
SELECT
CASE WHEN (su.user_id IS NULL OR CAST(ed.timestamp AS DECIMAL(20,0)) > CAST(su.subscrn_end_dt AS DECIMAL(20,0))) THEN 'UNSUBSCRIBED'
WHEN (su.user_id IS NOT NULL AND CAST(ed.timestamp AS DECIMAL(20,0)) <= CAST(su.subscrn_end_dt AS DECIMAL(20,0))) THEN 'SUBSCRIBED'
END AS user_type,
SUM(ABS(CAST(ed.end_ts AS DECIMAL(20,0))-CAST(ed.start_ts AS DECIMAL(20,0)))) AS duration
FROM enriched_data ed
LEFT OUTER JOIN subscribed_users su
ON ed.user_id=su.user_id
WHERE ed.status='pass'
AND ed.batchid='$batchId'
GROUP BY CASE WHEN (su.user_id IS NULL OR CAST(ed.timestamp AS DECIMAL(20,0)) > CAST(su.subscrn_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.subscrn_end_dt AS DECIMAL(20,0))) THEN 'SUBSCRIBED' END"""

|   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 = """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"""

|   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 = """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"""

|   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 = """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.subscrn_end_dt AS DECIMAL(20,0))))
GROUP BY ed.user_id
ORDER BY duration DESC
LIMIT 10"""

try {
    sqlContext.sql( sqlText = "SET hive.auto.convert.join=false")
    sqlContext.sql( sqlText = "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()
}
}
}

```

After executing the above script, we would get the below

```

18/07/24 22:00:50 INFO scheduler.OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator s
18/07/24 22:00:50 INFO spark.SparkContext: Successfully stopped SparkContext
18/07/24 22:00:50 INFO util.ShutdownHookManager: Shutdown hook called
18/07/24 22:00:50 INFO util.ShutdownHookManager: Deleting directory /tmp/spark-3c8ddd16-fb74-4688-adf2-0288d2e0da92
Running script for data analysis for different use cases... AnkithTest-Completed
USE CASES COMPLETE !!
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost music]$ 

```

And for each of the data analysis we have ingested the final output into the tables in HIVE.
Please find the new tables available in HIVE.

```

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: 0.06 seconds, Fetched: 11 row(s)
hive> 

```

```
hive> select * from connected_artists;
OK
A301      4          1
A302      3          1
A303      1          1
A304      1          1
Time taken: 0.341 seconds, Fetched: 4 row(s)
```

```
hive> select * from top_10_royalty_songs;
OK
S206      35231627      1
S208      5231627 1
S205      5231627 1
S201      2627294 1
S207      2627294 1
S200      2627294 1
S202      100000 1
Time taken: 0.389 seconds, Fetched: 7 row(s)
```

```
hive> select * from top_10_stations;
OK
ST400      2          2          1
ST404      1          1          1
ST403      1          1          1
ST410      1          1          1
ST407      1          1          1
ST413      1          1          1
Time taken: 0.435 seconds, Fetched: 6 row(s)
```

```
Time taken: 0.294 seconds, Fetched: 10 row(s)
hive> select * from top_10_unsubscribed_users;
OK
U111    28807006      1
U107    26202673      1
U119    26202673      1
U112    20000000      1
U113    20000000      1
U116    19900000      1
U117    10000000      1
U100    7835960  1
U114    5231627  1
U118    5231627  1
Time taken: 0.294 seconds, Fetched: 10 row(s)
```

```
hive> select * from usersBehaviour;
OK
UNSUBSCRIBED    174666154      1
SUBSCRIBED     81434300      1
Time taken: 0.424 seconds, Fetched: 2 row(s)
hive>
```

So we have successfully completed the data analysis for the given dataset. But now – we look forward to store the results. This is explained in the next section – we have implemented it by using the “export” concept of sqoop as mentioned in the script - **data_export.sh**.

4.5 Data Export

In **data_export.sh** we are going to export the data from the hive tables into mysql using **Sqoop** export. **data_export.sh**

```

#!/bin/bash
#This script is not working.
#Either change table to text or use STRING as type of partitioned column
batchid=`cat /home/acadgild/examples/music/logs/current-batch.txt`
LOGFILE=/home/acadgild/examples/music/logs/log_batch_$batchid
echo "Creating mysql tables if not present..." >> $LOGFILE
mysql -u "root" "-pRoot@123" < /home/acadgild/examples/music/create_schema.sql
echo "Running sqoop job for data export...AnkitTest" >> $LOGFILE
sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table top_10_stations --export-dir /user/hive/warehouse/project.db/top_10_stations/batchid=$batchid --input-fields-terminated-by ',' -m 1
sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table usersBehaviour --export-dir /user/hive/warehouse/project.db/usersBehaviour/batchid=$batchid --input-fields-terminated-by ',' -m 1
sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table connected_artists --export-dir /user/hive/warehouse/project.db/connected_artists/batchid=$batchid --input-fields-terminated-by ',' -m 1
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=$batchid --input-fields-terminated-by ',' -m 1
sqoop export --connect jdbc:mysql://localhost/project --username root --password Root@123 --table top_10_unsubscribed_users --export-dir /user/hive/warehouse/project.db/top_10_unsubscribed_users/batchid=$batchid --input-fields-terminated-by ',' -m 1

create_schema.sql
CREATE DATABASE IF NOT EXISTS project;

USE project;

CREATE TABLE IF NOT EXISTS top_10_stations
(
station_id VARCHAR(50),
total_distinct_songs_played INT,
distinct_user_count INT
);

CREATE TABLE IF NOT EXISTS usersBehaviour
(
user_type VARCHAR(50),
duration BIGINT
);

CREATE TABLE IF NOT EXISTS connected_artists
(
artist_id VARCHAR(50),
user_count INT
);

CREATE TABLE IF NOT EXISTS top_10_royalty_songs
(
song_id VARCHAR(50),
duration BIGINT
);

CREATE TABLE IF NOT EXISTS top_10_unsubscribed_users
(
user_id VARCHAR(50),
duration BIGINT
);

```

Now we can see the export using export sqoop is in-progress.

```

Let us export the final data.....
mysql: [Warning] Using a password on the command line interface can be insecure.
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 $HCAT_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 $ACCUMULO_HOME to the root of your Accumulo installation.
18/07/25 21:34:03 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6
18/07/25 21:34:03 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
18/07/25 21:34:03 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
18/07/25 21:34:03 INFO tool.CodeGenTool: Beginning code generation
Wed Jul 25 21:34:03 IST 2018 WARN: Establishing SSL connection without server's identity verification is not recommended. According to MySQL 5.5.45+, 5.6.26+ and 5.7.6+ 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=true and provide truststore for server certificate verification.
18/07/25 21:34:04 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `top_10_stations` AS t LIMIT 1
18/07/25 21:34:04 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `top_10_stations` AS t LIMIT 1
18/07/25 21:34:04 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /home/acadgild/install/hadoop/hadoop-2.6.5
Note: /tmp/sqoop-acadgild/compile/b03f92d636d323d60b5f4f15a9a82d55/top_10_stations.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
18/07/25 21:34:07 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-acadgild/compile/b03f92d636d323d60b5f4f15a9a82d55/top_10_stations.jar
18/07/25 21:34:07 INFO mapreduce.ExportJobBase: Beginning export of top_10_stations

Container
18/07/25 21:34:11 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1532530385941_0002
18/07/25 21:34:12 INFO impl.YarnClientImpl: Submitted application application_1532530385941_0002
18/07/25 21:34:12 INFO mapreduce.Job: The url to track the job: http://localhost:8088/proxy/application_1532530385941_0002
18/07/25 21:34:12 INFO mapreduce.Job: Running job: job_1532530385941_0002
18/07/25 21:34:21 INFO mapreduce.Job: Job job_1532530385941_0002 running in uber mode : false
18/07/25 21:34:21 INFO mapreduce.Job: map 0% reduce 0%
18/07/25 21:34:30 INFO mapreduce.Job: map 100% reduce 0%
18/07/25 21:34:30 INFO mapreduce.Job: Job job_1532530385941_0002 completed successfully
18/07/25 21:34:30 INFO mapreduce.Job: Counters: 30
    File System Counters
        FILE: Number of bytes read=0
        FILE: Number of bytes written=127642
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=276
        HDFS: Number of bytes written=0
        HDFS: Number of read operations=4
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=0
    Job Counters
        Launched map tasks=1
        Data-local map tasks=1
        Total time spent by all maps in occupied slots (ms)=5658
        Total time spent by all reduces in occupied slots (ms)=0
        Total time spent by all map tasks (ms)=5658
        Total vcore-milliseconds taken by all map tasks=5658
        Total megabyte-milliseconds taken by all map tasks=5793792
    Map-Reduce Framework
        Map input records=6
        Map output records=6
        Input split bytes=213
        Spilled Records=0
        Failed Shuffles=0
        Merged Map outputs=0
        GC time elapsed (ms)=63
        CPU time spent (ms)=1450

```

The sqoop export command exported the tables from the hive and it stored in the Mysql.

The data stored in the Mysql is shown in the successive screen shots –

“project” database

```
mysql> show databases;
+-----+
| Database           |
+-----+
| information_schema |
| metastore          |
| mysql              |
| oozie              |
| performance_schema |
| project            |
| simplidb          |
| sys                |
+-----+
8 rows in set (0.00 sec)
```

```
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 |
| top_10_royalty_songs |
| top_10_stations    |
| top_10_unsubscribed_users |
| users_beaviour     |
+-----+
5 rows in set (0.00 sec)
```

```
mysql> select * from connected_artists;
+-----+-----+
| artist_id | user_count |
+-----+-----+
| A301      |        4 |
| A302      |        3 |
| A303      |        1 |
| A304      |        1 |
+-----+-----+
4 rows in set (0.00 sec)
```

```
mysql> select * from top_10_royalty_songs;
+-----+-----+
| song_id | duration |
+-----+-----+
| S206    | 35231627 |
| S208    | 5231627  |
| S205    | 5231627  |
| S201    | 2627294  |
| S207    | 2627294  |
| S200    | 2627294  |
| S202    | 100000   |
+-----+-----+
7 rows in set (0.00 sec)
```

```
mysql> select * from top_10_stations;
+-----+-----+-----+
| station_id | total_distinct_songs_played | distinct_user_count |
+-----+-----+-----+
| ST400      |                2 |                  2 |
| ST404      |                1 |                  1 |
| ST403      |                1 |                  1 |
| ST410      |                1 |                  1 |
| ST407      |                1 |                  1 |
| ST413      |                1 |                  1 |
+-----+-----+-----+
6 rows in set (0.00 sec)
```

```
mysql> select * from top_10_unsubscribed_users;
+-----+-----+
| user_id | duration |
+-----+-----+
| U111    | 28807006 |
| U107    | 26202673 |
| U119    | 26202673 |
| U112    | 20000000 |
| U113    | 20000000 |
| U116    | 19900000 |
| U117    | 10000000 |
| U100    | 7835960  |
| U114    | 5231627  |
| U118    | 5231627  |
+-----+-----+
10 rows in set (0.00 sec)
```

```
mysql> select * from usersBehaviour;
+-----+-----+
| user_type | duration |
+-----+-----+
| UNSUBSCRIBED | 174666154 |
| SUBSCRIBED   | 81434300 |
+-----+-----+
2 rows in set (0.00 sec)
```

Now having a look at the log file generated.

```
[acadgild@localhost music]$ cat /home/acadgild/examples/music/logs/log_batch_1
Starting daemons
Creating LookUp Tables
Populating LookUp Tables
Creating hive tables on top of hbase tables for data enrichment and filtering...
Running script for Data Formatting...
Running script for data enrichment and filtering...
Copying valid and invalid records in local file system...
Deleting older valid and invalid records from local file system...
Creating mysql tables if not present...
Running sqoop job for data export...
```

4.6 Job Scheduling

For scheduling, the master file - music_project_master.sh, which holds all the scripts used in the Project, should run in a period of every 3 hours.

This scheduling can usually be done using the crontab. As of now, we are not going to schedule the file since this project was a self-learning activity where we had only one input batch files.

4.7 Problems faced during project installation and how it resolved

1. The frequent error was the metastore service was not connected due to which I had trouble on creation of hbase tables.
2. To execute the metastore service we have to make sure , mysql is up and running.
3. to execute the shell scripts we had to make sure those scripts have permission to execute, Or can be modified by using chmod command.