Big Data And Hadoop

Project 2.2 - USA Consumer Forum Data Analysis

1. **Executive Summary**
   **1.1. Project Overview**
   To develop the System to analyze the USA Consumer Forum Data.

   **1.2.Purpose And Scope of specification**
   The purpose of this project is to capture the data for analyzing the complaints
   registered by the consumers at the forum.

   **In Scope:**
   The following requirement will be addressed:
   a. Developing system to handle the data of the complaints registered at the forum
      and store the information in Hadoop Cluster (Flume).
   b. Analyze the data that could be used to get the insight of various businesses.
   c. Store the results in RDBMS (MySQL).

2. **Product/Service Description**
   **2.1.Assumptions**
   The data is accessible to the system to download it. The data is present in a CSV file.

   **2.2.Constraints**
   Describe any item that will constrain the design options, including
   a. This system may not be used for searching for now. But it will be used for
      analysis and saving the relevant information as of now.
   b. System will be using MySQL as the database.

3. **Problem Statement:**

The dataset is in csv format and contains the attributes pertaining to resolution of the
consumer complaints.

You need to copy the dataset into HDFS using Flume and the results of the problem
statements should be exported into RDBMS(Mysql) using sqoop.

The aim of this project is to analyze performance of various companies on aspects like:
i.      Write a pig script to find no of complaints which got timely response.
ii.     Write a pig script to find no of complaints where consumer forum forwarded the
        complaint same day they received to respective company.
iii.    Write a pig script to find list of companies toping in complaint chart (companies with
        maximum number of complaints).

iv.      Write a pig script to find no of complaints filed with product type has "Debt collection" for the year 2015.

Submit the screenshots of all the solutions with the source code.

## 4. Input dataset download link

Associated Data File is placed at the following location

*https://drive.google.com/file/d/0B1QaXx7tpw3SQTlnQ0MzVW5HajA/view?usp=sharing*

## 5. Dataset description

Below is the description of the data set

| Column heading | Index | Description |
|---|---|---|
| Date received | 0 | date on which consumer filed the complaint |
| Product | 1 | Type of the product |
| Sub-product | 2 | Sub product type |
| Issue | 3 | Issue faced by the consumer |
| Sub-issue | 4 | Any sub issues if exists |
| Consumer complaint narrative | 5 | Detailed description of complaint |
| Company public response | 6 | Company's public response to the complaint |
| Company | 7 | Name of the company |
| State | 8 | State from which consumer filed the complaint |
| ZIP code | 9 | Zip code |
| Submitted via | 10 | Channel from which complaint was submitted |
| Date sent to company | 11 | Date on which consumer forum forwarded the complaint to company |
| Company response to consumer | 12 | Company's response to the consumer |
| Timely response? | 13 | |
| Consumer disputed? | 14 | |
| Complaint ID | 15 | Unique complaint id |

This data is comma delimited.

## 6. Solution

### a. Download the data into local file system

The data file is downloaded from the link mentioned in section 4 and placed in the local file system at **/home/acadgild/Project2.2**. Name of the file is **Consumer_Complaints.csv.**

This file can be seen on the local system using **ls** command as follows:



Using **head** command, we get can view the content of this file as follows. Here we have displayed only first 50 lines of the file.



b.  **Start Hadoop cluster and history server**
    We first start Hadoop cluster using the command **start-all.sh** as follows:

Now, we start the history server using the command **mr-jobhistory-daemon.sh start historyserver** as follows:



The threads running after the start up are seen using **jps** command as follows:



c. **Load data into HDFS using Flume**

Now, we need to place this file from local file system to HDFS.

We use the **ls** command to check if the input file already existed in HDFS:
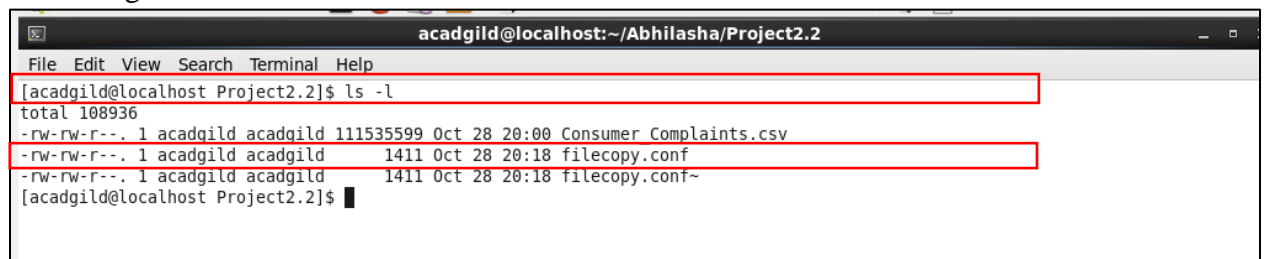
In the above screenshot, the file named **Consumer_Complaints.csv** is not listed. So, this file is not already present.

Flume agent is made of three parts
   a.  Source
   b.  Channel
   c.  Sink
In our use-case, source is the exec, sink in HDFS and the channel is memory channel.

The location of the configuration file to be used in flume is **/home/acadgild/Abhilasha/Project2.2**. Its name is **filecopy.conf**. This file can be seen using **ls** command as follows:

```
acadgild@localhost:~/Abhilasha/Project2.2
File  Edit  View  Search  Terminal  Help
[acadgild@localhost Project2.2]$ ls -l
total 108936
-rw-rw-r--. 1 acadgild acadgild 111535599 Oct 28 20:00 Consumer_Complaints.csv
-rw-rw-r--. 1 acadgild acadgild      1411 Oct 28 20:18 filecopy.conf
-rw-rw-r--. 1 acadgild acadgild      1411 Oct 28 20:18 filecopy.conf~
[acadgild@localhost Project2.2]$
```

P**roperties of source** defined in it are as follows:
i.   type = exec
ii.  command = hadoop dfs -put /home/acadgild/Abhilasha/Project2.2/Consumer_Complaints.csv /abhilasha/

P**roperties of the channel** defined in it are as follows:
i.   type = memory

**Properties of the sink** defined in it are as follows:
i.   type = hdfs
ii.  path = hdfs://localhost:9000/abhilasha/

The command used to put data into HDFS using flume is
*flume-ng agent –n agent1 –c conf –f /home/acadgild/Abhilasha/Project2.2/filecopy.conf*

In this command, we have mentioned the configuration file to be used in the execution of flume job.

After the successful execution of this command, we see the input file placed in HDFS using **ls** command.



We can also see it through the HDFS UI as follows:

**Task 1: Write a pig script to find no of complaints which got timely response**

**Answer:**

The fields from the data that we are going to use are *Complaint ID* and *Timely Response* flag.

The pig script that we are using to solve the problem statement has the following steps:

Step 1: *REGISTER '/usr/local/pig/lib/piggybank.jar';*

We are going to use *CSVExcelStorage* to read data from the csv file. In order to use this, which is present in *piggybank.jar,* we register this jar.

We mention the full qualified local path of the jar, which is */usr/local/pig/lib/piggybank.jar* in this case.

Step 2*: complaintDetails  = LOAD '/abhilasha/Consumer_Complaints.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',','NO_MULTILINE','UNIX','SKIP_ INPUT_HEADER') AS (dateRec:chararray, product:chararray, subProduct:chararray, issue:chararray, subIssue:chararray, complaintNarrative:chararray, companyPublicResponse:chararray, company:chararray, state:chararray, zipCode:chararray, submittedVia:chararray, dateSentToCompany:chararray, companyResponseToConsumer:chararray, timelyResponse:chararray, consumerDisputed:chararray, complaintId:chararray);*

This command is to specify the details of the input file, schema of the data if known, to be used to load the data.

Here, the file path given is */abhilasha/Consumer_Complaints.csv.* This is the location in HDFS. Rest is the schema of the input, as we have the schema beforehand.

Step 3: *complaintsWithTimelyResponse = FILTER complaintDetails BY timelyResponse == 'Yes';*

As we need records that were given timely response, we use *FILTER* to do so. The clause used to identify complaints with timely response is *timelyResponse == 'Yes'.*

Step 4: *distinctComplaints = GROUP complaintsWithTimelyResponse BY complaintId;*

This command is used to get rid of the redundant records. It is found that the data in the input file had a pinch of duplicity. Hence, this is the measure taken to avoid duplicate complaint IDs in the data. So, we group the data by complaint ID.

Step 5: *complaintIds = FOREACH distinctComplaints GENERATE group AS complaintID;*

Of all the fields in the data, now only complaint ID remains of use to us. Hence, we extract only the complaint ID from the records in the previous step.

Step 6: *groupForCount = GROUP complaintIds ALL;*

This step is to aid us in getting the total count of complaints. This step will get all the complaint IDs in a single group.

Step 7: *count = FOREACH groupForCount GENERATE COUNT(complaintIds);*

This is used to get the count of complaints grouped in the previous step.

Step 8: *STORE count INTO '/abhilasha/Project2.2.Task1';*

This is to store the data back into HDFS. The destination path given is */abhilasha/Project2.2.Task1*.

All these commands are put together in a file named *Task1* stored at */home/acadgild/Abhilasha/Project2.2.*

This script file is executed as follows:

The command used is *pig Task1.pig*

The command indicates that the script will be run not locally but will use HDFS to read and write data.

The underlying map-reduce job can be seen successfully completed in the job history server as follows:



Details of job execution can be seen below:

The output path mentioned in the script was */abhilasha/Project2.2.Task1* in HDFS. We can see the output folder that got created as a result of the job execution as follows:
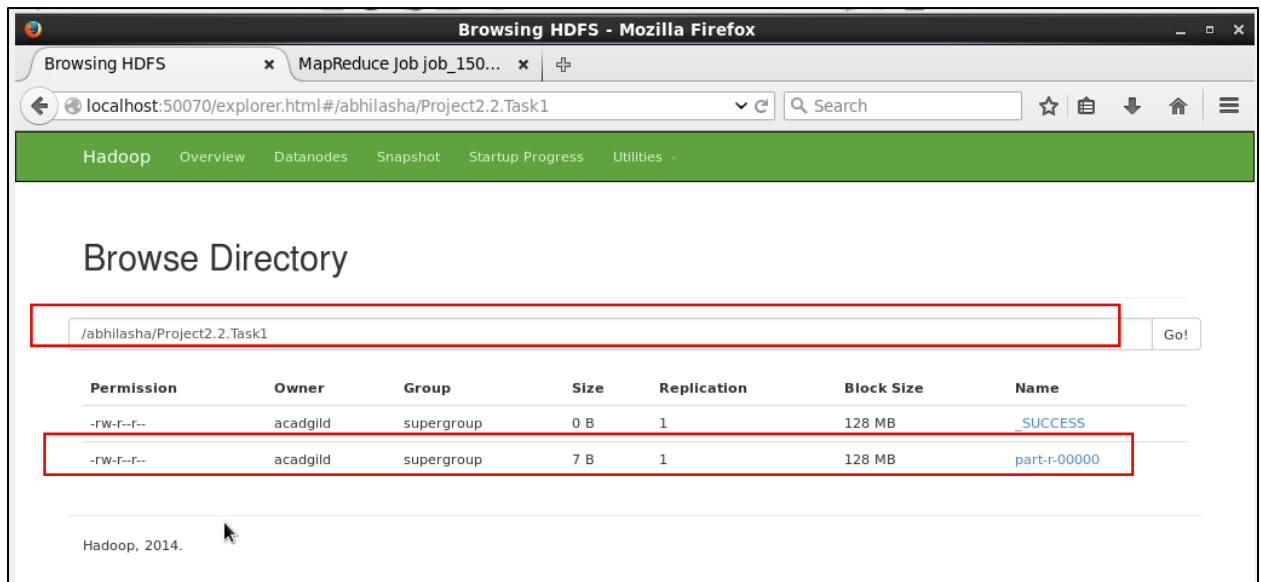


This folder contains the files mentioned in the screen shot below. Of these *part-r-00000* file contains the output of the pig script.
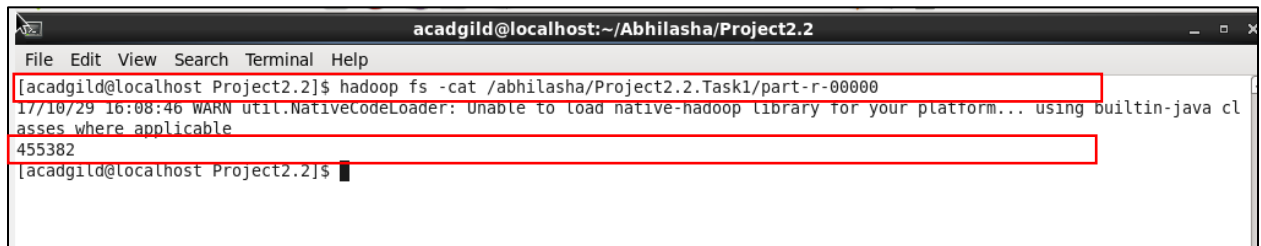
These files can also be viewed from HDFS UI as shown below:



The content of this output file can be seen using the **cat** command as follows:



**Note: To upload the output file in GitHub, we have renamed the file to** *part-r-000001.*

Now that the output is placed in HDFS, we need to store the output in MySQL. To do so, we perform the following steps

Step 1: Start MySQL using the command
*sudo service mysqld start*



Step 2: To start the command line interface of MySQL, we use the command
*mysql -u root*

Step 3: We create a new database named Project22 using the command
*create database Project22;*



Step 4: We change the database from default to Project22, we use the command

*use Project22;*



Step 5: We now create table named *Task1* to store the output of Task 1 in it using the command
*create table Task1*
*(*
      *CountOfComplaints int*
*);*
Here, the column of the table is *CountOfComplaints.*

Step 6: The created table can be listed using the command

*Show tables;*



Step 6: Next step is to export the data from HDFS and store in into the table created in previous step.

*sqoop export --connect jdbc:mysql://localhost/Project22 --username 'root' -P --table 'Task1' --export-dir '/abhilasha/Project2.2.Task1/part-r-00000' -m 1 --columns CountOfComplaints;*

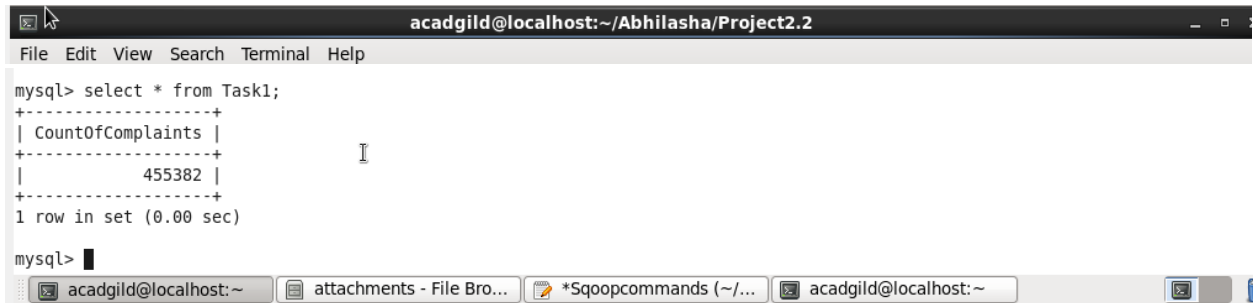Here, the parameters mentioned in the above command are:

i.     Database : jdbc:mysql://localhost/Project22
ii.    Table Name : Task1
iii.   Path where exported file is placed : /abhilasha/Project2.2.Task1/part-r-00000
iv.    Column to populate : CountOfComplaints

The command execution is shown below:



Step 7: After the successful placement of data in MySQL, we can see the content of the table in the database populated, using *select * from Task1;* as follows:

```
acadgild@localhost:~/Abhilasha/Project2.2
File  Edit  View  Search  Terminal  Help
mysql> select * from Task1;
+-------------------+
| CountOfComplaints |
+-------------------+
|            455382 |
+-------------------+
1 row in set (0.00 sec)

mysql>
```
acadgild@localhost:~    attachments - File Bro...   *Sqoopcommands (~/...   acadgild@localhost:~

This shows that the data is placed successfully in database.

---

**Task 2: Write a pig script to find no of complaints where consumer forum forwarded the complaint same day they received to respective company**

**Answer:**

The fields from the data that we are going to use are *Date Received*, *Date sent to company* and *Complaint ID*.

The pig script that we are using to solve the problem statement has the following steps:

Step 1: *REGISTER '/usr/local/pig/lib/piggybank.jar';*

We are going to use *CSVExcelStorage* to read data from the csv file. In order to use this, which is present in *piggybank.jar,* we register this jar.

We mention the full qualified local path of the jar, which is */usr/local/pig/lib/piggybank.jar* in this case.

Step 2*: complaintDetails  = LOAD '/abhilasha/Consumer_Complaints.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',','NO_MULTILINE','UNIX','SKIP_ INPUT_HEADER') AS (dateRec:chararray, product:chararray, subProduct:chararray, issue:chararray, subIssue:chararray, complaintNarrative:chararray, companyPublicResponse:chararray, company:chararray, state:chararray, zipCode:chararray, submittedVia:chararray, dateSentToCompany:chararray, companyResponseToConsumer:chararray, timelyResponse:chararray, consumerDisputed:chararray, complaintId:chararray);*

This command is to specify the details of the input file, schema of the data if known, to be used to load the data.

Here, the file path given is */abhilasha/Consumer_Complaints.csv.* This is the location in HDFS. Rest is the schema of the input, as we have the schema beforehand.

Step 3: *complaintsRequired = FILTER complaintDetails BY dateRec == dateSentToCompany;*

As we need records where consumer forum forwarded the complaint same day they received to respective company, we use the clause *dateRec == dateSentToCompany* to filter the records.

Step 4: *distinctComplaints = GROUP complaintsRequired BY complaintId;*

This command is used to get rid of the redundant records. It is found that the data in the input file had a pinch of duplicity. Hence, this is the measure taken to avoid duplicate complaint IDs in the data. So, we group the data by complaint ID.

Step 5: *complaintIds = FOREACH distinctComplaints GENERATE group AS complaintID;*

Of all the fields in the data, now only complaint ID remains of use to us. Hence, we extract only the complaint ID from the records in the previous step.

Step 6: *groupForCount = GROUP complaintIds ALL;*

This step is to aid us in getting the total count of complaints. This step will get all the complaint IDs in a single group.

Step 7: *count = FOREACH groupForCount GENERATE COUNT(complaintIds);*

This is used to get the count of complaints grouped in the previous step.

Step 8: *STORE count INTO '/abhilasha/Project2.2.Task2';*

This is to store the data back into HDFS. The destination path given is */abhilasha/Project2.2.Task2.*

All these commands are put together in a file named *Task2* stored at */home/acadgild/Abhilasha/Project2.2.*

This script file is executed as follows:

The command used is *pig Task2.pig*

The command indicates that the script will be run not locally but will use HDFS to read and write data.



The underlying map-reduce job can be seen successfully completed in the job history server as follows:



Details of job execution can be seen below:

The output path mentioned in the script was */abhilasha/Project2.2.Task2* in HDFS. We can see the output folder that got created as a result of the job execution as follows:



This folder contains the files mentioned in the screen shot below. Of these *part-r-00000* file contains the output of the pig script.

These files can also be viewed from HDFS UI as shown below:



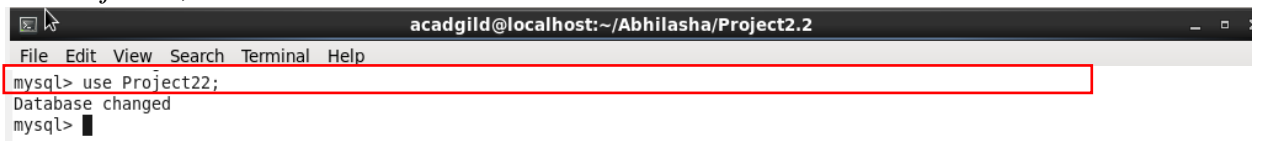The content of this output file can be seen using the **cat** command as follows:



**Note: To upload the output file in GitHub, we have renamed the file to *part-r-000002.***

Now that the output is placed in HDFS, we need to store the output in MySQL. To do so, we perform the following steps:

We have already seen steps to start MySQL and create database *Project22*. So not mentioning those steps again.

Step 1: As seen previously, we have created the database named *Project22*. We change the database from default to Project22, we use the command

*use Project22;*

Step 5: We now create table named *Task2* to store the output of Task 2 in it using the command
*create table Task2*
*(*
      *CountOfComplaints int*
*);*
Here, the column of the table is *CountOfComplaints.*



Step 6: The created table can be listed using the command

*Show tables;*



Step 6: Next step is to export the data from HDFS and store in into the table created in previous step.
*sqoop export --connect jdbc:mysql://localhost/Project22 --username 'root' -P --table 'Task2' --export-dir '/abhilasha/Project2.2.Task2/part-r-00000' -m 1 --columns CountOfComplaints;*
Here, the parameters mentioned in the above command are:
i.        Database : jdbc:mysql://localhost/Project22
ii.      Table Name : Task2
iii.     Path where exported file is placed : /abhilasha/Project2.2.Task2/part-r-00000
iv.     Column to populate : CountOfComplaints

The command execution is shown below:

Step 7: After the successful placement of data in MySQL, we can see the content of the table in the database populated, using *select \* from Task2;* as follows:



This shows that the data is placed successfully in database.

---

**Task 3: Write a pig script to find list of companies toping in complaint chart (companies with maximum number of complaints).**

**Answer:**

The fields from the data that we are going to use are *company* and *Complaint ID*.

The pig script that we are using to solve the problem statement has the following steps:

Step 1: *REGISTER '/usr/local/pig/lib/piggybank.jar';*

We are going to use *CSVExcelStorage* to read data from the csv file. In order to use this, which is present in *piggybank.jar,* we register this jar.

We mention the full qualified local path of the jar, which is
*/usr/local/pig/lib/piggybank.jar* in this case.

Step 2*: complaintDetails  = LOAD '/abhilasha/Consumer_Complaints.csv' USING
org.apache.pig.piggybank.storage.CSVExcelStorage(',','NO_MULTILINE','UNIX','SKIP_
INPUT_HEADER') AS (dateRec:chararray, product:chararray, subProduct:chararray,
issue:chararray, subIssue:chararray, complaintNarrative:chararray,
companyPublicResponse:chararray, company:chararray, state:chararray,
zipCode:chararray, submittedVia:chararray, dateSentToCompany:chararray,
companyResponseToConsumer:chararray, timelyResponse:chararray,
consumerDisputed:chararray, complaintId:chararray);*

This command is to specify the details of the input file, schema of the data if known, to
be used to load the data.

Here, the file path given is */abhilasha/Consumer_Complaints.csv.* This is the location in
HDFS. Rest is the schema of the input, as we have the schema beforehand.

Step 3: *companywiseComplaints = GROUP complaintDetails BY company;*

Here, we group the records by company.

Step 4: *companywiseCount = FOREACH companywiseComplaints GENERATE group
AS companyName, COUNT(complaintDetails.complaintId) AS countOfComplaints;*

After the grouping of records is done by company, we find the count of records per
group. We have used the alias *countOfComplaints* to store the count. This is done for
every company in  *companywiseComplaints* and hence, used *FOREACH.*

Step 5: *orderByCount = ORDER companywiseCount BY countOfComplaints DESC;*

We now order the records in descending order of the count of complaints.

Step 6: *companyWithMaxComplaints = LIMIT orderByCount 1;*

We need the company with highest count of complaints only. Hence, we limit the data to
1 record.

Step 7: *STORE companyWithMaxComplaints INTO '/abhilasha/Project2.2.Task3' USING PigStorage('|');*

This is to store the data back into HDFS. The destination path given is */abhilasha/Project2.2.Task3*. Also, the delimiter used to separate the fields in the output file is '|'.

All these commands are put together in a file named *Task3* stored at */home/acadgild/Abhilasha/Project2.2*.

This script file is executed as follows:

The command used is *pig Task3.pig*

The command indicates that the script will be run not locally but will use HDFS to read and write data.



The underlying map-reduce job can be seen successfully completed in the job history server as follows:



Details of job execution can be seen below:

The output path mentioned in the script was */abhilasha/Project2.2.Task3* in HDFS. We can see the output folder that got created as a result of the job execution as follows:

This folder contains the files mentioned in the screen shot below. Of these *part-r-00000* file contains the output of the pig script.



These files can also be viewed from HDFS UI as shown below:



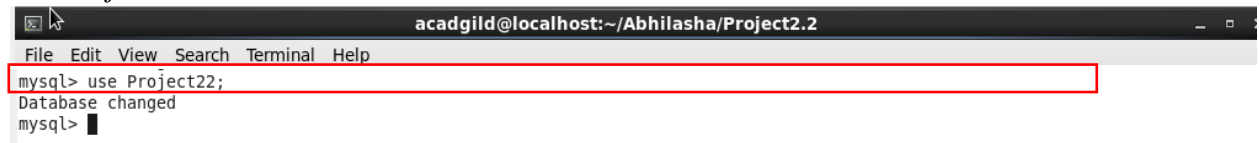The content of this output file can be seen using the **cat** command as follows:



**Note: To upload the output file in GitHub, we have renamed the file to *part-r-000003*.**

Now that the output is placed in HDFS, we need to store the output in MySQL. To do so, we perform the following steps:

We have already seen steps to start MySQL and create database *Project22*. So not mentioning those steps again.

Step 1: As seen previously, we have created the database named *Project22*. We change the database from default to Project22, we use the command

*use Project22;*



Step 5: We now create table named *Task2* to store the output of Task 2 in it using the command
*create table Task3*
*(*

        *Company varchar(2000),*
        *CountOfComplaints int*

*);*
Here, the columns of the table are *Company* and *CountOfComplaints*.



Step 6: The created table can be listed using the command
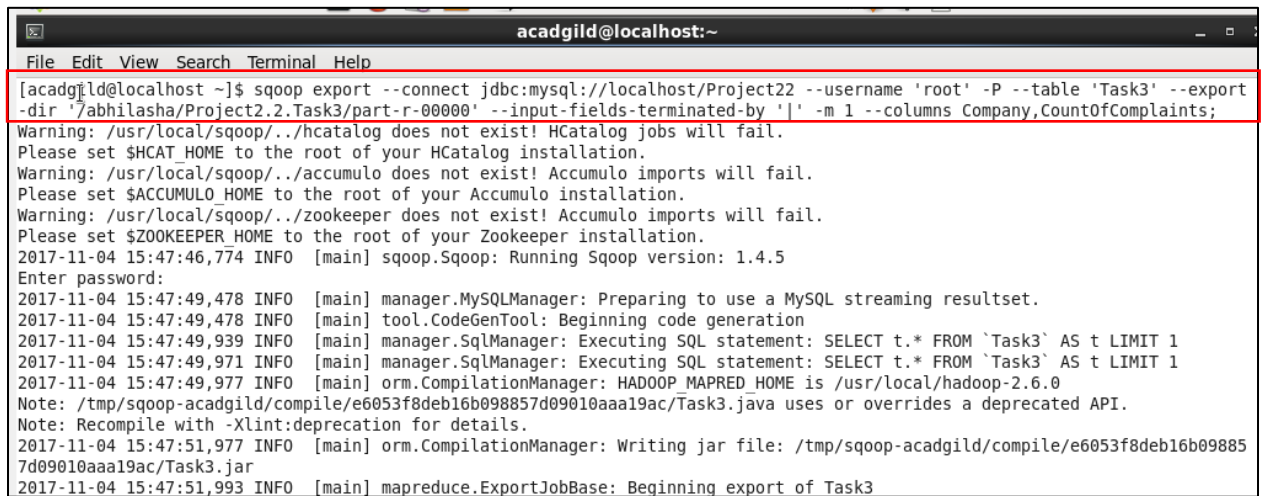
*Show tables;*



Step 6: Next step is to export the data from HDFS and store in into the table created in previous step.

*sqoop export --connect jdbc:mysql://localhost/Project22 --username 'root' -P --table 'Task3' --export-dir '/abhilasha/Project2.2.Task3/part-r-00000' –input-field-terminated-by '|' -m 1 --columns Columns,CountOfComplaints;*

Here, the parameters mentioned in the above command are:

i.      Database : jdbc:mysql://localhost/Project22
ii.     Table Name : Task2
iii.    Path where exported file is placed : /abhilasha/Project2.2.Task2/part-r-00000
iv.     Field separator : |
v.      Columns to populate : Company, CountOfComplaints

The command execution is shown below:



Step 7: After the successful placement of data in MySQL, we can see the content of the table in the database populated, using *select * from Task3;* as follows:



This shows that the data is placed successfully in database.

---

**Task 4: Write a pig script to find no of complaints filed with product type has "Debt collection" for the year 2015.**

**Answer:**

The fields from the data that we are going to use are *Date Received*, *Product to company* and *Complaint ID*.

The pig script that we are using to solve the problem statement has the following steps:

Step 1: *REGISTER '/usr/local/pig/lib/piggybank.jar';*

We are going to use *CSVExcelStorage* to read data from the csv file. In order to use this, which is present in *piggybank.jar,* we register this jar.

We mention the full qualified local path of the jar, which is */usr/local/pig/lib/piggybank.jar* in this case.

Step 2*: complaintDetails = LOAD '/abhilasha/Consumer_Complaints.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',','NO_MULTILINE','UNIX','SKIP_ INPUT_HEADER') AS (dateRec:chararray, product:chararray, subProduct:chararray, issue:chararray, subIssue:chararray, complaintNarrative:chararray, companyPublicResponse:chararray, company:chararray, state:chararray, zipCode:chararray, submittedVia:chararray, dateSentToCompany:chararray, companyResponseToConsumer:chararray, timelyResponse:chararray, consumerDisputed:chararray, complaintId:chararray);*

This command is to specify the details of the input file, schema of the data if known, to be used to load the data.

Here, the file path given is */abhilasha/Consumer_Complaints.csv*. This is the location in HDFS. Rest is the schema of the input, as we have the schema beforehand.

Step 3: *requiredRecords = FILTER complaintDetails BY (product == 'Debt collection' AND dateRec MATCHES '.*2015.*');*

As we need records where complaints filed with product type has " Debt collection" for the year 2015, we use the clause *product == 'Debt collection' AND dateRec MATCHES '.*2015.*',* to filter the records. Here, *MATCHES* is used to perform pattern matching to get records where date of registration of complaint is 2015.

Step 4: *distinctComplaints = GROUP requiredRecords BY complaintId;*

This command is used to get rid of the redundant records. It is found that the data in the input file had a pinch of duplicity. Hence, this is the measure taken to avoid duplicate complaint IDs in the data. So, we group the data by complaint ID.

Step 5: *complaintIds = FOREACH distinctComplaints GENERATE group AS complaintID;*

Of all the fields in the data, now only complaint ID remains of use to us. Hence, we extract only the complaint ID from the records in the previous step.

Step 6: *groupForCount = GROUP complaintIds ALL;*

This step is to aid us in getting the total count of complaints. This step will get all the complaint IDs in a single group.

Step 7: *count = FOREACH groupForCount GENERATE COUNT(complaintIds);*

This is used to get the count of complaints grouped in the previous step.

Step 8: *STORE count INTO '/abhilasha/Project2.2.Task4';*

This is to store the data back into HDFS. The destination path given is */abhilasha/Project2.2.Task4.*

All these commands are put together in a file named *Task4* stored at */home/acadgild/Abhilasha/Project2.2.*

This script file is executed as follows:

The command used is *pig Task4.pig*

The command indicates that the script will be run not locally but will use HDFS to read and write data.

The underlying map-reduce job can be seen successfully completed in the job history server as follows:



Details of job execution can be seen below:

The output path mentioned in the script was */abhilasha/Project2.2.Task4* in HDFS. We can see the output folder that got created as a result of the job execution as follows:
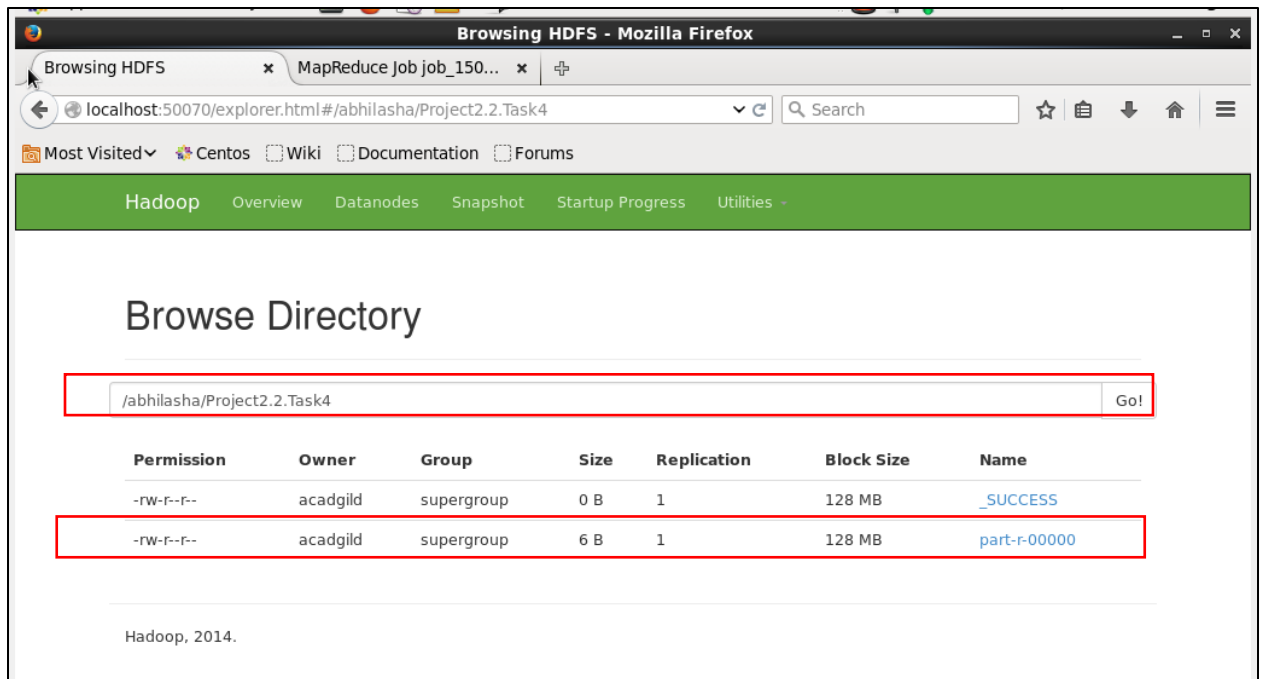


This folder contains the files mentioned in the screen shot below. Of these *part-r-00000* file contains the output of the pig script.
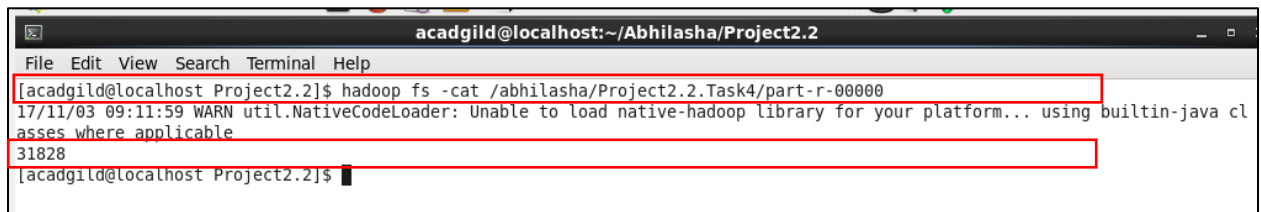


These files can also be viewed from HDFS UI as shown below:

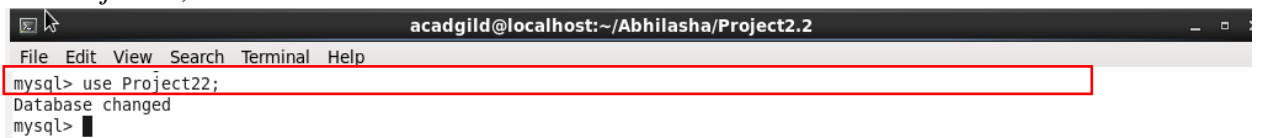The content of this output file can be seen using the **cat** command as follows:



**Note: To upload the output file in GitHub, we have renamed the file to** *part-r-000004.*

Now that the output is placed in HDFS, we need to store the output in MySQL. To do so, we perform the following steps:

We have already seen steps to start MySQL and create database *Project22.* So not mentioning those steps again.

Step 1: As seen previously, we have created the database named *Project22.* We change the database from default to Project22, we use the command

*use Project22;*

Step 5: We now create table named *Task4* to store the output of Task 4 in it using the command
*create table Task4*
*(*
*CountOfComplaints int*
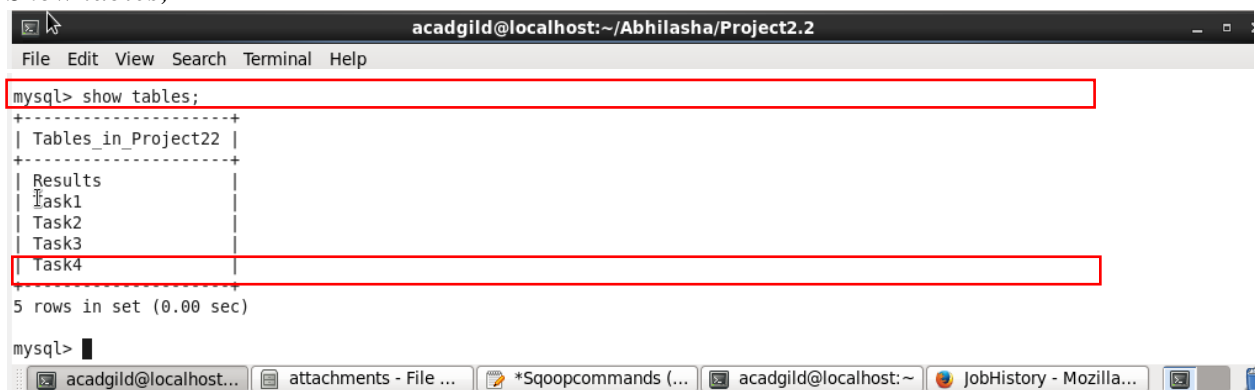*);*
Here, the column of the table is *CountOfComplaints.*



Step 6: The created table can be listed using the command

*Show tables;*



Step 6: Next step is to export the data from HDFS and store in into the table created in previous step.
*sqoop export --connect jdbc:mysql://localhost/Project22 --username 'root' -P --table*
*'Task4' --export-dir '/abhilasha/Project2.2.Task4/part-r-00000' -m 1 --columns*
*CountOfComplaints;*
Here, the parameters mentioned in the above command are:
v.       Database : jdbc:mysql://localhost/Project22
vi.      Table Name : Task4
vii.     Path where exported file is placed : /abhilasha/Project2.2.Task4/part-r-00000
viii.    Column to populate : CountOfComplaints

The command execution is shown below:

Step 7: After the successful placement of data in MySQL, we can see the content of the table in the database populated, using *select * from Task4;* as follows:



This shows that the data is placed successfully in database.