# Assessment 4: Managing data with Hive and Pig in Hadoop

## MA5831 – Advanced Data Processing and Analysis using SAS

## 13848336 Nikki Fitzherbert

## Technical Preparation and Setup

A couple of preparatory steps were required before any querying of the consumer complaints data could be performed. These are outlined briefly in this and the next section.

The following code extract from the MRemoteNG session shows how it was verified that the required data file (*consumer_complaints.txt*) did in fact exist, how a new folder called *consumer* was created in the *DIHPS* database in the Hadoop Distributed File System (HDFS) and how a copy of the data file was moved to the new folder.

```
Using username "student".
[student@server3 ~]$ pwd
/home/student
[student@server3 ~]$ ls -R /home/student/DIHPS/data
/home/student/DIHPS/data:
census_2010.csv  census_2012.csv  census_2014.csv  consumer_complaints.txt  wordcount.pig
census_2011.csv  census_2013.csv  census.csv        numbers.txt
[student@server3 ~]$ hdfs dfs -ls /user/student/DIHPS
Found 6 items
drwxr-xr-x   - student hive          0 2020-09-21 05:45 /user/student/DIHPS/census_data
drwxr-xr-x   - student hive          0 2020-09-19 19:03 /user/student/DIHPS/data
drwxr-xr-x   - student hive          0 2020-09-21 07:40 /user/student/DIHPS/output
drwxr-xr-x   - student hive          0 2020-09-21 05:51 /user/student/DIHPS/population_census
drwxr-xr-x   - student hive          0 2020-09-21 06:08 /user/student/DIHPS/population_census2
drwxr-xr-x   - student hive          0 2020-09-19 19:16 /user/student/DIHPS/test_table
[student@server3 ~]$ hdfs dfs -mkdir -p /user/student/DIHPS/consumer
[student@server3 ~]$ hdfs dfs -put /home/student/DIHPS/data/consumer_complaints.txt /user/student/DIHPS/consumer
[student@server3 ~]$ hdfs dfs -ls /user/student/DIHPS
Found 7 items
drwxr-xr-x   - student hive          0 2020-09-21 05:45 /user/student/DIHPS/census_data
drwxr-xr-x   - student hive          0 2020-09-22 08:35 /user/student/DIHPS/consumer
drwxr-xr-x   - student hive          0 2020-09-19 19:03 /user/student/DIHPS/data
drwxr-xr-x   - student hive          0 2020-09-21 07:40 /user/student/DIHPS/output
drwxr-xr-x   - student hive          0 2020-09-21 05:51 /user/student/DIHPS/population_census
drwxr-xr-x   - student hive          0 2020-09-21 06:08 /user/student/DIHPS/population_census2
drwxr-xr-x   - student hive          0 2020-09-19 19:16 /user/student/DIHPS/test_table
[student@server3 ~]$ hdfs dfs -ls -R /user/student/DIHPS/consumer
-rw-r--r--   1 student hive   59878874 2020-09-22 08:35 /user/student/DIHPS/consumer/consumer_complaints.txt
```
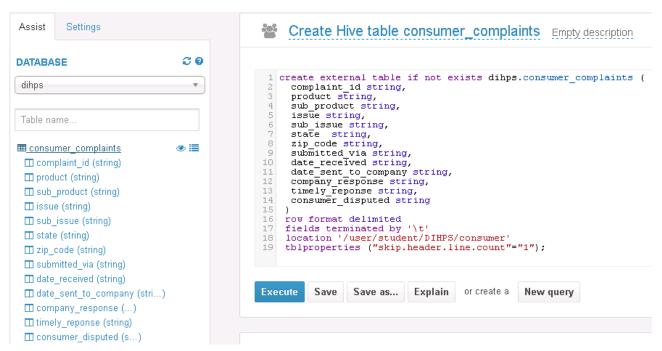
This extract from Hue shows that a copy of the *consumer complaints* data was successfully transferred into the HDFS.

**Data Processing Tasks**

The final data preparation step was to create a Hive table schema appropriate to support any subsequent query operations using HiveQL. Therefore, an external table called *consumer_complaints* was created using the 'create table' command in Hue, with the location of the underlying data indicated by the 'location' property, and the format of the stored data indicated by the 'row format' and 'fields terminated by' properties. The two advantages of creating an external table rather than a managed table were it was more flexible in that other tools could also query the data, and the source data remained intact and untouched.

Note that for simplicity, all fields were read in as strings and the 'tblproperties' option was used to skip the first row of the underlying data file as it contained the field names.



The following extract shows that the Hive table was successfully created and populated with the expected consumer complaints data.



| | complaint_id | product | sub_product | issue | sub_issue | state | zip_code | submitted_via | date_received | date_sent_to_company | company_response | tir |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1431865 | Consumer loan | Vehicle loan | Managing the loan or lease | | NJ | 8736 | Web | 6/22/2015 | 6/22/2015 | In progress | Yes |
| 1 | 1431374 | Debt collection | Medical | Disclosure verification of debt | Not given enough info to verify debt | WI | 54140 | Web | 6/22/2015 | 6/22/2015 | Closed with explanation | Yes |
| 2 | 1431251 | Mortgage | Conventional fixed mortgage | Loan modification,collection,foreclosure | | MO | 63368 | Web | 6/22/2015 | 6/22/2015 | In progress | Yes |
| 3 | 1431743 | Debt collection | Medical | Cont'd attempts collect debt not owed | Debt is not mine | WA | 98055 | Web | 6/22/2015 | 6/22/2015 | Closed with explanation | Yes |
| 4 | 1432678 | Debt collection | Medical | Cont'd attempts collect debt not owed | Debt was paid | TX | 75104 | Web | 6/22/2015 | 6/22/2015 | Closed with explanation | Yes |
| 5 | 1432104 | Debt collection | Other (phone, health club, etc.) | Cont'd attempts collect debt not owed | Debt was paid | CA | 95423 | Web | 6/22/2015 | 6/22/2015 | In progress | Yes |
| 6 | 1431998 | Debt collection | Other (phone, health club, etc.) | Communication tactics | Frequent or repeated calls | TX | 75048 | Web | 6/22/2015 | 6/22/2015 | Closed with explanation | Yes |
| 7 | 1432207 | Debt collection | Medical | Disclosure verification of debt | Right to dispute notice not received | TX | 75125 | Web | 6/22/2015 | 6/22/2015 | Closed with explanation | Yes |
| 8 | 1432202 | Consumer loan | Vehicle loan | Problems when you are unable to pay | | TN | 37174 | Web | 6/22/2015 | 6/22/2015 | In progress | Yes |
| 9 | 1432334 | Debt collection | Other (phone, health club, etc.) | Disclosure verification of debt | Right to dispute notice not received | FL | 33830 | Web | 6/22/2015 | 6/22/2015 | Closed with non-monetary relief | Yes |
| 10 | 1431188 | Debt collection | Payday loan | Communication tactics | Threatened to take legal action | NE | 68134 | Web | 6/22/2015 | 6/22/2015 | Closed with explanation | Yes |
| 11 | 1431358 | Money transfers | Domestic (US) money transfer | Money was not available when promised | | AL | 35235 | Phone | 6/22/2015 | 6/22/2015 | In progress | Yes |

# Reporting and Coding Tasks – Hive Queries

## Question 1  Determine the 10 states with the maximum number of complaints

```
1  select state, count(complaint_id) as num_complaints
2  from dihps.consumer_complaints
3  group by state
4  order by num_complaints desc
5  limit 10;
```

Execute   Save   Save as...   Explain   or create a   New query

## Question 2  Determine how many complaints are associated with the 'Medical' sub-product offering

```
1  select count(complaint_id) as num_complaints
2  from dihps.consumer_complaints
3  where sub_product = 'Medical';
```

Execute   Save   Save as...   Explain   or create a   New query

## Question 3  Determine five ZIP codes with the smallest number of complaints

```
1  select zip_code, count(complaint_id) as num_complaints
2  from dihps.consumer_complaints
3  where zip_code > 500
4  group by zip_code
5  order by num_complaints asc
6  limit 5;
```

Execute   Save   Save as...   Explain   or create a   New query

## Question 3b  Determine total number of valid ZIP codes with one complaint

```
1  select count(x.num_complaints) as num_complaints from (
2    select zip_code, count(complaint_id) as num_complaints
3    from dihps.consumer_complaints
4    where zip_code > 500
5    group by zip_code
6    order by num_complaints asc
7  ) as x
8  where x.num_complaints = 1;
```

### Question 4 — Determine how many complaints, grouped by product and state, contain the word 'fraud' in the issue description

```sql
select product, state, count(complaint_id) as num_complaints
from dihps.consumer_complaints
where lower(issue) like '%fraud%'
group by product, state;
```

Execute   Save   Save as...   Explain   or create a   New query

### Question 5 — Create a table summarising the total number of complaints by product, state and submitted_via

```sql
create table if not exists dihps.summary as
select product, state, submitted_via, count(*) as num_complaints
from dihps.consumer_complaints
group by product, state, submitted_via;
```

**Reporting and Coding Tasks – Pig Latin scripts**

## Question 6 - Number of complaints submitted via web vs other methods

```
1  -- load consumer_complaints text file from HDFS location using a tab as a delimiter
2  T = LOAD '/user/student/DIHPS/consumer/consumer_complaints.txt' using PigStorage('\t') AS
3  (
4  -- define scheme to be used to read the text file
5  complaint_id:chararray,
6  product:chararray,
7  sub_product:chararray,
8  issue:chararray,
9  sub_issue:chararray,
10 state:chararray,
11 zip_code:chararray,
12 submitted_via:chararray,
13 date_received:chararray,
14 date_sent_to_compant:chararray,
15 company_response:chararray,
16 timely_response:chararray,
17 consumer_disputed:chararray
18 );
19
20 --filter rows of file to include only those where lower-case submitted_via equals 'web'
21 T_WEB = FILTER T BY LOWER(submitted_via) == 'web';
22 -- filter rows of ffile to include only those where lower-case submitted_via does not equal 'web'
23 T_NWEB = FILTER T BY LOWER(submitted_via) != 'web';
24
25 -- store the results in two separate tables in the HDFS DIHPS output folder
26 -- web results table
27 STORE T_WEB INTO '/user/student/DIHPS/output/web_results';
28 -- non-web results table
29 STORE T_NWEB INTO '/user/student/DIHPS/output/other_results';
30
31 -- verify the number of rows written to each table
32 -- group the results together
33 T_WEB_GRP = GROUP T_WEB ALL;
34 T_NWEB_GRP = GROUP T_NWEB ALL;
35
36 -- count the number of rows in the files
37 T_WEB_COUNT = FOREACH T_WEB_GRP GENERATE COUNT(T_WEB) AS num_complaints;
38 T_NWEB_COUNT = FOREACH T_NWEB_GRP GENERATE COUNT(T_NWEB) AS num_complaints;
39
40 -- store the results in two separate tables in the HDFS DIHPS output folder
41 -- web results table count
42 STORE T_WEB_COUNT INTO '/user/student/DIHPS/output/web_results_count';
43 -- non-web results table count
44 STORE T_NWEB_COUNT INTO '/user/student/DIHPS/output/other_results_count';
```

## Question 7 - Top 10 states with the maximum number of complaints

```
1  -- load consumer_complaints text file from HDFS location using a tab as a delimiter
2  T = LOAD '/user/student/DIHPS/consumer/consumer_complaints.txt' using PigStorage('\t') AS
3  (
4  -- define schema to be used to read the text file
5  complaint_id:chararray,
6  product:chararray,
7  sub_product:chararray,
8  issue:chararray,
9  sub_issue:chararray,
10 state:chararray,
11 zip_code:chararray,
12 submitted_via:chararray,
13 date_received:chararray,
14 date_sent_to_company:chararray,
15 company_response:chararray,
16 timely_response:chararray,
17 consumer_disputed:chararray
18 );
19
20 -- group rows of file by state field
21 TO = GROUP T BY state;
22 -- count the number of rows of data in each group
23 T_GRP = FOREACH TO GENERATE group, COUNT(T) AS state_count;
24 -- order the groups by the state counts in descending order
25 T_GRP2 = ORDER T_GRP BY state_count DESC;
26 -- limit the number of rows of output to 10
27 T_LIM = LIMIT T_GRP2 10;
28
29 -- store the results in a table called 'max_complaints' in the HDFS DIHPS folder
30 STORE T_LIM INTO '/user/student/DIHPS/output/max_complaints';
```

## Question 8 - Complaints associated with the 'Medical' sub-product offering

```
1  -- load consumer_complaints text file from HDFS location using a tab as a delimiter
2  T = LOAD '/user/student/DIHPS/consumer/consumer_complaints.txt' using PigStorage('\t') AS
3  (
4  -- define schema to be used to read the text file
5  complaint_id:chararray,
6  product:chararray,
7  sub_product:chararray,
8  issue:chararray,
9  sub_issue:chararray,
10 state:chararray,
11 zip_code:chararray,
12 submitted_via:chararray,
13 date_received:chararray,
14 date_sent_to_company:chararray,
15 company_response:chararray,
16 timely_response:chararray,
17 consumer_disputed:chararray
18 );
19
20 -- filter rows of file to include only those where sub-product = 'Medical'
21 TO = FILTER T BY sub_product == 'Medical';
22
23 -- store the results in a table called 'medical_complaints_list' in the HDFS DIHPS output folder
24 STORE TO INTO '/user/student/DIHPS/output/medical_complaints_list';
25
26 -- group the results together
27 T_GRP = GROUP TO ALL;
28
29 -- count the number of rows in the file
30 T_CNT = FOREACH T_GRP GENERATE COUNT(TO) AS num_complaints;
31
32 -- store the result in a table called 'medical_complaints_total' in the HDFS DIHPS output folder
33 STORE T_CNT INTO '/user/student/DIHPS/output/medical_complaints_total';
```

## Question 9 - Five zip codes with the least number of complaints

```
 1  -- load consumer_complaints text file from HDFS location using a tab as a delimiter
 2  T = LOAD '/user/student/DIHPS/consumer/consumer_complaints.txt' using PigStorage('\t') AS
 3  (
 4  -- define schema to be used to read the text file
 5  complaint_id:chararray,
 6  product:chararray,
 7  sub_product:chararray,
 8  issue:chararray,
 9  sub_issue:chararray,
10  state:chararray,
11  zip_code:chararray,
12  submitted_via:chararray,
13  date_received:chararray,
14  date_sent_to_company:chararray,
15  company_response:chararray,
16  timely_response:chararray,
17  consumer_disputed:chararray
18  );
19
20  -- group rows of file by zip code field
21  TO = GROUP T BY zip_code;
22  -- count the number of rows of data in each group
23  T_GRP = FOREACH TO GENERATE group, COUNT(T) AS zip_count;
24  -- order the groups by the zip code counts in ascending order
25  T_GRP2 = ORDER T_GRP BY zip_count ASC;
26  -- limit the number of rows of output to 5
27  T_LIM = LIMIT T_GRP2 5;
28
29  -- store the results in a table called 'least_complaints' in the HDFS DIHPS folder
30  STORE T_LIM INTO '/user/student/DIHPS/output/least_complaints';
```

## Question 10 - Complaints associated with identity theft

```
 1  -- load consumer_complaints text file from HDFS location using a tab as a delimiter
 2  T = LOAD '/user/student/DIHPS/consumer/consumer_complaints.txt' using PigStorage('\t') AS
 3  (
 4  -- define schema to be used to read the text file
 5  complaint_id:chararray,
 6  product:chararray,
 7  sub_product:chararray,
 8  issue:chararray,
 9  sub_issue:chararray,
10  state:chararray,
11  zip_code:chararray,
12  submitted_via:chararray,
13  date_received:chararray,
14  date_sent_to_company:chararray,
15  company_response:chararray,
16  timely_response:chararray,
17  consumer_disputed:chararray
18  );
19
20  -- filter rows of file to include only those where the issue description contains the word 'identity theft'
21  TO = FILTER T BY issue MATCHES '.*Identity theft.*';
22  -- group the result by the product and state fields
23  T_GRP = GROUP TO BY (product, state);
24
25  -- store the result in a table called 'id_theft_complaints' in the HDFS DIHPS folder
26  STORE T_GRP INTO '/user/student/DIHPS/output/id_theft_complaints';
```

```
1  -- load consumer_complaints text file from HDFS location using a tab as a delimiter
2  T = LOAD '/user/student/DIHPS/consumer/consumer_complaints.txt' using PigStorage('\t') AS
3  (
4  -- define schema to be used to read the text file
5  complaint_id:chararray,
6  product:chararray,
7  sub_product:chararray,
8  issue:chararray,
9  sub_issue:chararray,
10 state:chararray,
11 zip_code:chararray,
12 submitted_via:chararray,
13 date_received:chararray,
14 date_sent_to_company:chararray,
15 company_response:chararray,
16 timely_response:chararray,
17 consumer_disputed:chararray
18 );
19
20 -- group the rows by the product, sub-product and state fields
21 T_GRP = GROUP T BY (product, sub_product, state);
22
23 -- flatten each group and count the number of rows
24 T_SUMMARY = FOREACH T_GRP GENERATE FLATTEN(group), COUNT(T) AS num_complaints;
25
26 -- store the result in a table called 'complaint_summary' in the HDFS DIHPS folder
27 STORE T_SUMMARY INTO '/user/student/DIHPS/output/complaint_summary';
```