

****Problem Statement 2****

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

Step 1: Start pig in mapreduce mode.

\$ mr-jobhistory-daemon.sh start historyserver

\$ pig

```
[acadgild@localhost project2.1]$ mr-jobhistory-daemon.sh start historyserver
starting historyserver, logging to /usr/local/hadoop-2.6.0/logs/mapred-acadgild-historyserver-localhost.localdomain.out
```

```
[acadgild@localhost project2.1]$ pig
2017-08-25 17:08:53,786 INFO [main] pig.ExecTypeProvider: Trying ExecType : LOCAL
2017-08-25 17:08:53,789 INFO [main] pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
2017-08-25 17:08:53,789 INFO [main] pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2017-08-25 17:08:53,904 [main] INFO org.apache.pig.Main - Apache Pig version 0.14.0 (r1640057) compiled Nov 16 2014, 18:02:05
2017-08-25 17:08:53,904 [main] INFO org.apache.pig.Main - Logging error messages to: /home/acadgild/project2.1/pig_1503661133904.log
2017-08-25 17:08:53,951 [main] INFO org.apache.pig.impl.util.Utils - Default bootup file /home/acadgild/.pigbootup not found
2017-08-25 17:08:54,378 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2017-08-25 17:08:54,378 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2017-08-25 17:08:54,378 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://localhost:9000
2017-08-25 17:08:54,383 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.used.genericoptionsparser is deprecated. Instead, use mapreduce.client.genericoptionsparser.used
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/local/hbase/lib/slf4j-log4j12-1.6.4.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop-2.6.0/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.
2017-08-25 17:08:54,689 [main] WARN org.apache.hadoop.util.NativeCodeLoader - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2017-08-25 17:08:55,270 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt>
```

Step 2: register piggybank.jar:

REGISTER piggybank.jar;

DEFINE CSVExcelStorage org.apache.pig.piggybank.storage.CSVExcelStorage;

```
grunt> REGISTER piggybank.jar;
grunt>
grunt> DEFINE CSVExcelStorage org.apache.pig.piggybank.storage.CSVExcelStorage;
grunt>
```

CSVExcelStorage will be used to load CSV file into pig relation.

Step 3: Load file into relation:

A = LOAD '/flume_sink2/' USING CSVExcelStorage('','NO_MULTILINE','UNIX','SKIP_INPUT_HEADER');*

describe A;

```
grunt> A = LOAD '/flume_sink2/*' USING CSVExcelStorage('','NO_MULTILINE','UNIX','SKIP_INPUT_HEADER');
2017-08-26 16:49:09,750 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-08-26 16:49:09,750 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-08-26 16:49:09,751 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> describe A;
Schema for A unknown.
```

Relation 'A' has been loaded with CSV file along with below details:

Separator Character: “,”

Multi line treatment of the record: *NO_MULTILINE*

Line Break Type: UNIX

Header of CSV: SKIP_INPUT_HEADER = this will ensure that header of the CSV will not be loaded into relation.

Step 4: Select the required columns from dataset:

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.

```
B = FOREACH A GENERATE (int)$15 AS complainID,
```

```
(chararray)$0 AS datereceived,
```

```
(chararray)$11 AS dateforwarded,
```

```
(chararray)$7 AS companyName;
```

```
grunt> B = FOREACH A GENERATE (int)$15 AS complainID,  
>> (chararray)$0 AS datereceived,  
>> (chararray)$11 AS dateforwarded,  
>> (chararray)$7 AS companyName;  
grunt> █
```

Step 5: Filter the relation based on datereceived==dateforwarded

```
C = FILTER B BY datereceived==dateforwarded;
```

```
grunt> C = FILTER B BY datereceived==dateforwarded;  
grunt> █
```

Step 6: Calculate the count of complains by grouping on companyName:

```
D = GROUP C BY companyName;
```

```
E = FOREACH D GENERATE group as companyName, COUNT(C.complainID) as complainCount;
```

```
grunt> D = GROUP C BY companyName;
grunt> E = FOREACH D GENERATE group as companyName, COUNT(C.complainID) as complainCount;
grunt> █
```

Step 7: Store the result in '/user/acadgild/project/USAConsumer/ProblemStatement2' HDFS directory:

STORE E INTO '/user/acadgild/project/USAConsumer/ProblemStatement2';

```
2017-08-26 18:31:09,357 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2017-08-26 18:31:09,357 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.2.0 0.14.0 acadgild 2017-08-26 18:30:43 2017-08-26 18:31:09 GROUP_BY,FILTER

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_1503739494588_0027 1 1 8 8 8 4 4 4 A,B,C,D,E GROUP_BY,COMBINER /user/acadgild/project/USAConsumer/ProblemStatement2,

Input(s):
Successfully read 0 records from: "/flume_sink2/"

Output(s):
Successfully stored 0 records in: "/user/acadgild/project/USAConsumer/ProblemStatement2"

Counters:
Total records written : 0
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1503739494588_0027

2017-08-26 18:31:09,359 [main] INFO org.apache.hadoop.yarn.client.RMPProxy - Connecting to ResourceManager at /0.0.0.0:8032
2017-08-26 18:31:09,362 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2017-08-26 18:31:09,415 [main] INFO org.apache.hadoop.yarn.client.RMPProxy - Connecting to ResourceManager at /0.0.0.0:8032
2017-08-26 18:31:09,420 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2017-08-26 18:31:09,464 [main] INFO org.apache.hadoop.yarn.client.RMPProxy - Connecting to ResourceManager at /0.0.0.0:8032
2017-08-26 18:31:09,466 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2017-08-26 18:31:09,509 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Unable to retrieve job to compute warning aggregation
2017-08-26 18:31:09,509 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt> █
```

Step 8: Check the content in HDFS directory:

\$ hadoop fs -ls /user/acadgild/project/USAConsumer/ProblemStatement2/

*\$ hadoop fs -cat /user/acadgild/project/USAConsumer/ProblemStatement2/**

```
[acadgild@localhost project2.2]$ hadoop fs -ls /user/acadgild/project/USAConsumer/ProblemStatement2/
17/08/26 18:32:47 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 acadgild supergroup 0 2017-08-26 18:31 /user/acadgild/project/USAConsumer/ProblemStatement2/_SUCCESS
-rw-r--r-- 1 acadgild supergroup 62450 2017-08-26 18:31 /user/acadgild/project/USAConsumer/ProblemStatement2/part-r-00000
[acadgild@localhost project2.2]$ hadoop fs -cat /user/acadgild/project/USAConsumer/ProblemStatement2/*
17/08/26 18:33:04 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
FTC 2
OCC 1
Amex 1751
HSBC 1253
MEFA 19
VHDA 6
WLCC 1
Xoom 12
Ocwen 6251
Sarima 9
Sigur 1
Avante 15
MOHELA 32
Netnet 17
PayPal 459
DCQ LLC 3
Equifax 11968
Innovis 29
JMD,LLC 1
NVR Inc 15
Navient 2979
Regions 407
Seterus 824
TD Bank 1051
Barclays 818
CIT Bank 4
CMM, LLC 1
Ceannate 19
Citibank 6933
Comerica 105
ConServe 83
DLC, LLC 3
Discover 1588
```

PIG output have been stored successfully with the data separated by TAB (\t) which shows companyName and Count of complaints which are forwarded on the same day received by consumer forum.

PIG script is save with name ProblemStatement2.PIG which can be run using below command:

\$ pig <filepath>/ProblemStatement2.PIG

SQOOP: Export the results to mysql.

Step 1: start mysql/services:

\$ sudo service mysqld status

\$ sudo service mysqld start

\$ sudo service mysqld status

```
[acadgild@localhost project2.1]$ sudo service mysqld status
[sudo] password for acadgild:
mysqld is stopped
[acadgild@localhost project2.1]$ sudo service mysqld start
Starting mysqld: [ OK ]
[acadgild@localhost project2.1]$ sudo service mysqld status
mysqld (pid 9463) is running...
[acadgild@localhost project2.1]$
```

\$ mysql -u root

```
[acadgild@localhost project2.1]$ mysql -u root
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 2
Server version: 5.1.73 Source distribution

Copyright (c) 2000, 2013, 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>
```

Above command launches mysql with user root.

Step2: create table ForwardSameDaywith column companyName & complaintCount:

use db1;

show tables;

create table ForwardSameDay

(

companyName varchar(75),

complaintCount int(6)

);

describe ForwardSameDay;

*select * from ForwardSameDay;*

```
mysql> use db1;
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_db1 |
+-----+
| TimelyResponse |
| customer        |
| statewiseBPL80  |
| statewiseBPLacheived |
+-----+
4 rows in set (0.00 sec)

mysql> create table ForwardSameDay
-> (
->   companyName varchar(75),
->   complaintCount int(6)
-> );
Query OK, 0 rows affected (0.01 sec)

mysql> describe ForwardSameDay;
+-----+-----+-----+-----+-----+-----+
| Field          | Type          | Null | Key | Default | Extra |
+-----+-----+-----+-----+-----+-----+
| companyName    | varchar(75)   | YES  |     | NULL    |       |
| complaintCount | int(6)        | YES  |     | NULL    |       |
+-----+-----+-----+-----+-----+-----+
2 rows in set (0.00 sec)

mysql> select * from ForwardSameDay;
Empty set (0.00 sec)

mysql>
```

Step 3: run sqoop export command to get data from output directory of the pig job to mysql table.

```
sqoop export --connect jdbc:mysql://localhost/db1 \

--username 'root' -P --table 'ForwardSameDay' \

--export-dir '/user/acadgild/project/USAConsumer/ProblemStatement2/' \

--input-fields-terminated-by '\t' \

-m 1
```

```
[acadgild@localhost project2.2]$ sqoop export --connect jdbc:mysql://localhost/db1 \
> --username 'root' -P --table 'ForwardSameDay' \
> --export-dir '/user/acadgild/project/USAConsumer/ProblemStatement2/' \
> --input-fields-terminated-by '\t' \
> -m 1
Warning: /usr/local/sqoop/./hcatalog does not exist! HCatalog jobs will fail.
Please set $HCAT_HOME to the root of your HCatalog installation.
Warning: /usr/local/sqoop/./accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
Warning: /usr/local/sqoop/./zookeeper does not exist! Accumulo imports will fail.
Please set $ZOOKEEPER_HOME to the root of your Zookeeper installation.
2017-08-26 18:39:44,353 INFO [main] sqoop.Sqoop: Running Sqoop version: 1.4.5
Enter password:
2017-08-26 18:39:52,792 INFO [main] manager.MySQLManager: Preparing to use a MySQL streaming resultset.
2017-08-26 18:39:52,793 INFO [main] tool.CodeGenTool: Beginning code generation
2017-08-26 18:39:53,099 INFO [main] manager.SqlManager: Executing SQL statement: SELECT t.* FROM `ForwardSameDay` AS t LIMIT 1
2017-08-26 18:39:53,131 INFO [main] manager.SqlManager: Executing SQL statement: SELECT t.* FROM `ForwardSameDay` AS t LIMIT 1
2017-08-26 18:39:53,138 INFO [main] orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/local/hadoop-2.6.0
2017-08-26 18:39:57,157 INFO [main] mapreduce.JobSubmitter: Submitting tokens for job: job_1503739494588_0028
2017-08-26 18:39:57,564 INFO [main] impl.YarnClientImpl: Submitted application application_1503739494588_0028 to ResourceManager at /0.0.0.0:8032
2017-08-26 18:39:57,658 INFO [main] mapreduce.Job: The url to track the job: http://localhost:8088/proxy/application_1503739494588_0028/
2017-08-26 18:39:57,659 INFO [main] mapreduce.Job: Running job: job_1503739494588_0028
2017-08-26 18:40:05,854 INFO [main] mapreduce.Job: Job job_1503739494588_0028 running in uber mode : false
2017-08-26 18:40:05,857 INFO [main] mapreduce.Job: map 0% reduce 0%
2017-08-26 18:40:11,909 INFO [main] mapreduce.Job: map 100% reduce 0%
2017-08-26 18:40:11,916 INFO [main] mapreduce.Job: Job job_1503739494588_0028 completed successfully
```

Sqoop command completed successfully.

Step 4: check table in mysql:

*select * from ForwardSameDay;*

```
mysql> select * from ForwardSameDay;
```

companyName	complaintCount
FTC	2
OCC	1
Amex	1751
HSBC	1253
MEFA	19
VHDA	6
WLCC	1
Xoom	12
Ocwen	6251
Sarma	9
Sigue	1
Avante	15
MOHELA	32
Nelnet	17
PayPal	459
DCQ LLC	3
Equifax	11968
Innovis	29
JMD, LLC	1
NVR Inc	15
Navient	2979
Regions	407
Seterus	824
TD Bank	1051
Barclays	818
CIT Bank	4
CMM, LLC	1
Ceannate	19
Citibank	6933
Comerica	105
ConServe	83
DLC, LLC	3
Discover	1588

Step 5: Verify if all data have been exported from HDFS to MySQL:

Check number of lines in the HDFS file directory:

\$ hadoop fs -cat /user/acadgild/project/USAConsumer/ProblemStatement2/ | wc -l*

```
[acadgild@localhost project2.2]$ hadoop fs -cat /user/acadgild/project/USAConsumer/ProblemStatement2/* | wc -l
17/08/26 18:43:12 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin
2147
[acadgild@localhost project2.2]$
```

Check the count of the Table *ForwardSameDay* in mysql:

select count() from ForwardSameDay;*

```
mysql> select count(*) from ForwardSameDay;
+-----+
| count(*) |
+-----+
|      2147 |
+-----+
1 row in set (0.00 sec)

mysql>
```

As compared above all the data has been exported from HDFS to mysql using Sqoop.