# Spam Detection System in Hadoop

## 1) Cloud Infrastructure Setup

    a.   Create a Hadoop cluster using GCP (Google Cloud Platform) Dataproc.

I have chosen my cloud platform as GCP and I have created a cluster on it using Dataproc Compute Engine. The cluster is set to have 1 master and 3 worker nodes with the region selected as us-central1 and the zone as us-central1-f.
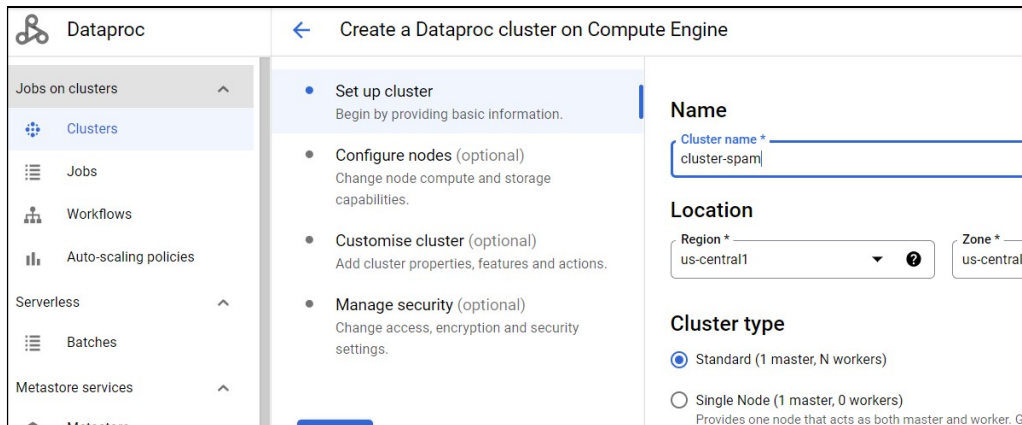


**Figure 1: Dataproc with 1master, N-worker**

To configure my nodes, I selected the machine from the family of General Purpose of N1 series and machine type as 'n1-standard-4 (4 vCPU, 15 GB memory)'. I have taken the same configuration for my worker nodes as well.
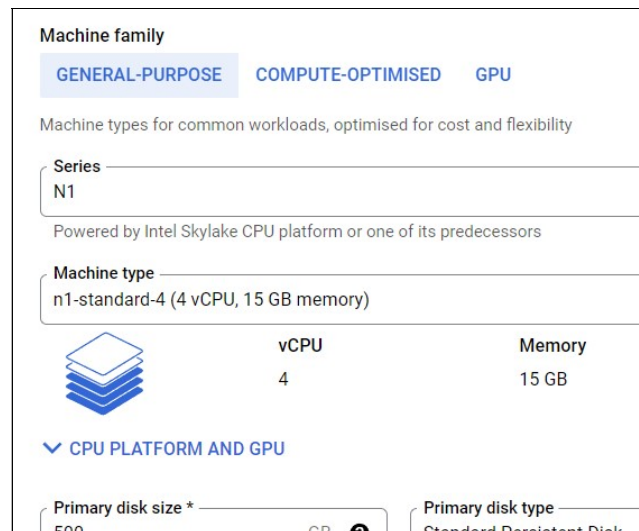


**Figure 2: Master node**

**Figure 3: Worker Node**

After the cluster has been created, we can verify the same from the Dataproc VM instances of our cluster.



**Figure 4: Cluster with 1 master, 3 worker**

b.  Environment check at the cluster.

In the dataproc cluster, Hadoop, pig, hive, and java was successfully installed. We have also created a directory of my user name under the '/user/'   with the command: hdfs dfs -mkdir /<user_name>.

We have also checked the version of the same to verify whether the installations have been properly done.



**Figure 5: Version check for Hadoop, Java, Pig, Hive**

I have checked whether the cluster has been set up properly by using the '~$ env' command.



**Figure 6: Environment check**

I have also provided access to professor Manoj on my cloud system, such that my project can by verified.



**Figure 7: Access control**

## 2) Choosing the dataset.

    a. Data from Kaggle.com

Link: https://www.kaggle.com/datasets/naveedhn/amazon-product-review-spam-and-non-spam

I have chosen a fairly large and complex dataset from Kaggle.com about the spam and ham reviews that different user on Amazon provides under the 'Cell phone and accessories' category.

**Figure 8: Dataset preview**

The dataset is of 1.87 Gb in size and of type .JSON and has 12 attributes which makes it more complex, and interesting to apply analytics on.

### b. Loading the data in the GCP bucket

I have loaded the data in the GCP bucket with the gsutil path as

"*gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0f a/cluster-spam-m/Cell_Phones_and_Accessories.json*"



**Figure 9: Dataset uploaded in Google Bucket**

## 3) Cleaning the data

The data has been cleaned has been performed with the following commands in Apache Pig.

```
>>cleanedtable1= FOREACH tableframe GENERATE id AS id, reviewerID AS reviewerID, reviewerName AS reviewerName,
REPLACE(reviewText,'([^a-zA-Z\\s]+)','') AS reviewText, class AS class;
```

```
>> cleanedtable1_5= FOREACH tableframe GENERATE id AS id, reviewerID AS reviewerID, reviewerName AS reviewerName,
REPLACE(reviewText, '[\r\n]+',' ') AS reviewText, class AS class;

>>cleanedtable2= FOREACH cleanedtable1_5 GENERATE id AS id, reviewerID AS reviewerID, reviewerName AS reviewerName,
LOWER(reviewText) AS reviewText, class AS class;

>>cleanedtable3 = FILTER cleanedtable2 BY reviewerName != 'Amazon Customer' AND reviewerName != 'No Name';

>>cleanedtable4 = FILTER cleanedtable3 BY reviewerName is not null AND reviewerID is not null;

>>cleanedtable5 = FILTER cleanedtable4 by reviewText is not null;

>>cleanedtable = FILTER cleanedtable5 by class==0 OR class==1;

>>STORE cleanedtable INTO
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/cleanedtable.json' USING JsonStorage();
```

In the first cleanedtable1 I used Regular Expression (Regex) to remove the special characters and numbers present in the review text. Also, I have restricted the unnecessary attributes from the original table, which made the new data more concise and useful for analysis.

In the cleanedtable1_5 I took the help of Regular Expression again to replace the new lines with spaces.

In cleanedtable2 I have loaded the review text by lowercasing the entire review text such that my analysis will be irrespective of the cases on the same word.

I have verified in cleanedtable3 that there are a fairly large amount of reviewers as 'Amazon Customer' and 'No Name', perhaps because there have been comments from Guest Accounts, so I am removing them to perform a better analysis.

In cleanedtable4 and cleanedtable5 I have removed all the rows where the reviewer name, id and text are null.

Finally since the 'class' is the attribute which classifies if the data is spam or ham, I have taken only those rows where the data has been already classified.

At the last I have stored the cleaned table in the cloud bucket such that when we go further on the analysis we will not need to perform the cleaning over and over again, we can directly load the data in the pig using the JsonLoader();



**Figure 10: Cleaned Table uploaded in Google Bucket**

## 4) Ham or Spam using Apache Pig

### a. Differentiate Ham and Spam Dataset

```
>>filter_spam = FILTER cleanedtable BY class == 1;
>>spamgroup= GROUP filter_spam ALL;
>>spamcount= FOREACH spamgroup GENERATE COUNT (filter_spam.reviewerID);
>>dump spamcount;

>>filter_ham = FILTER cleanedtable BY class == 0;
>>hamgroup= GROUP filter_ham ALL;
>>hamcount= FOREACH hamgroup GENERATE COUNT (filter_ham.reviewerID);
>>dump hamcount;


>>STORE filter_spam INTO
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/Spam_dataset.json' USING JsonStorage();

>>STORE filter_ham INTO
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/Ham_dataset.json' USING JsonStorage();
```
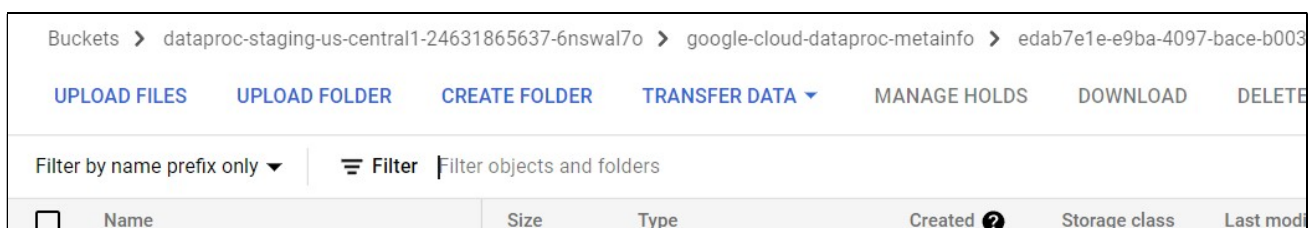
I have split the spam and the ham dataset in Apache pig using the 'class' attribute present in the original dataset, which represents whether a data is a spam or ham. After the split, I stored them separately in two variables filter_spam and filter_ham variable. To check whether the split has been performed correctly I have grouped the dataset and with foreach loop I have taken the count of the 'reviewerID', which is unique.



**Figure 11: Total count of Spam Data**



**Figure 12: Total count of Ham Data**

After that, I stored the 2 different datasets in the google cloud bucket, such that the data remains ready even after the instances are closed.

**Figure 13: Ham and Spam Dataset stored in Google Bucket**

b. Top 10 Spam Account

I have written the following code in Apache pig to find the top 10 Spam Account.

```
>>spamtest = LOAD
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/Spam_dataset.json' USING JsonLoader('id:(oid:chararray), reviewerID:chararray, reviewerName:chararray,
reviewText:chararray, summary:chararray, class:int');

>>spamReviewerGroup = GROUP spamtest by reviewerID;
>>spamReviewerCount = FOREACH spamReviewerGroup GENERATE FLATTEN(group) as reviewerID, COUNT($1) as count;
>>spamReviewerOrdered= ORDER spamReviewerCount by count DESC;
>>top10SpamReviewer= LIMIT spamReviewerOrdered 10;
>>SpamTotalData1= JOIN top10SpamReviewer by reviewerID, spamtest by reviewerID;
>>SpamNameAndId= FOREACH SpamTotalData1 GENERATE top10SpamReviewer::reviewerID as reviewerID, spamtest::reviewerName
as reviewerName, top10SpamReviewer::count as count;
Dump SpamNameAndId;
>>SpamNameAndId1= DISTINCT SpamNameAndId;
>>Dump SpamNameAndId1;

>>STORE SpamNameAndId1 INTO
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/Top10SpamAccount.json' USING JsonStorage();
```

At first, I loaded the spam data from the google cloud bucket, which we had previously saved. Furthermore, I created a new table frame by grouping the spam table with the reviewerId and adding the count of the rows in another column mentioned as 'count'. I have used 'FLATTEN()' operartor to un-nest the tuples which got created when I grouped spamtest by the reviewerId.
Following that I have sorted the table by descending order of the count because my aim is to find the Top 10 spam accounts and stored them in 'spamReviewerOrdered'.
Now, the challenge is the table frame, spamReviewerOrdered does not have the name of the reviewer, it only possesses the reviewerId and the count. So, in order to overcome this, I have joined the table frame top10SpamReviewer and spamtest with reviewerId as the pivot.
After the join is successfully done, with the FOREACH loop I have iterated through the table and picked only those attributes which are important to me i.e Reviewer ID, Reviewer Name and Count.
In the variable SpamNameAndId1, we have taken only the distinct values and printed in the console.

**Figure 14: Top 10 Spam Accounts**

a. Top 10 Ham Account

I have performed the steps with the same logic which I had previously performed while calculating the Top 10 Spam accounts.

The code is as follows.

```
>>hamtest = LOAD
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-ham-m/ham_dataset.json' USING JsonLoader('id:(oid:chararray), reviewerID:chararray, reviewerName:chararray,
reviewText:chararray, summary:chararray, class:int');

>>hamReviewerGroup = GROUP hamtest by reviewerID;

>>hamReviewerCount = FOREACH hamReviewerGroup GENERATE FLATTEN(group) as reviewerID, COUNT($1) as count;

>>hamReviewerOrdered= ORDER hamReviewerCount by count DESC;

>>top10hamReviewer= LIMIT hamReviewerOrdered 10;

>>hamTotalData1= JOIN top10hamReviewer by reviewerID, hamtest by reviewerID;

>>hamNameAndId= FOREACH hamTotalData1 GENERATE top10hamReviewer::reviewerID as reviewerID, hamtest::reviewerName as
reviewerName, top10hamReviewer::count as count;

>>hamNameAndId1= DISTINCT hamNameAndId;


>>STORE hamNameAndId1 INTO
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/Top10hamAccount.json' USING JsonStorage();
```
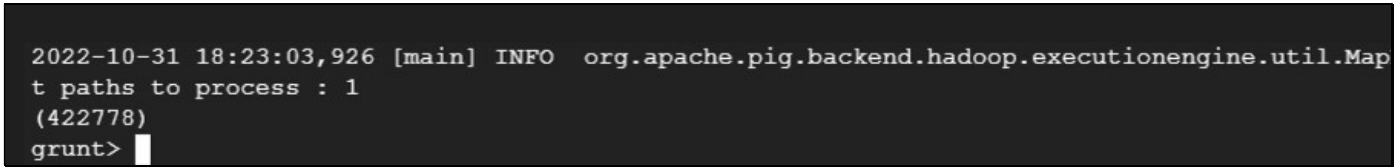
**Figure 15: Top 10 Ham Accounts**

I have saved both the top 10 Spam and Ham accounts in the google cloud bucket using JsonLoader()



**Figure 16: Top 10 Ham and Spam Accounts stored in Google Cloud Bucket**

## 5) TF-IDF using Map Reduce

### a. Top 10 Spam keywords for top 10 spam accounts.

The statistical technique known as Term Frequency - Inverse Document Frequency (TF-IDF) is frequently employed in information retrieval and natural language processing. It gauges a term's significance within a document in relation to a group of papers (i.e., relative to a corpus).

Formula-

- Tf=(Total number of times the term appears in a document/total number of terms in the document)

- Idf= Log(number of the document in the corpus/number of documents in the corpus containing the term)

- TfIdf= Tf*Idf

#### i. Data Preparation

In the data preparation step, we have taken the help of Apache Pig to load the spam data and find the top 10 spam reviewer (the logic has been previously explained when we calculated the top 10 spam reviewers).

In the variable SpamTotalData1 we have joined top10SpamReviewer and spamtest by reviewerId. After that from the resultant variable we have iterated the through the rows and only picked the attributes I require for the TF-IDF process i.e the reviewer name, reviewer id and the text and saved it the cloud storage.

So, the data which we possess right now has all the text comments by the top the spam accounts.

```
>>spamtest = LOAD
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/Spam_dataset.json' USING JsonLoader('id:(oid:chararray), reviewerID:chararray, reviewerName:chararray,
reviewText:chararray, summary:chararray, class:int');
```

```
>>spamReviewerGroup = GROUP spamtest by reviewerID;
>>spamReviewerCount = FOREACH spamReviewerGroup GENERATE FLATTEN(group) as reviewerID, COUNT($1) as count;
>>spamReviewerOrdered= ORDER spamReviewerCount by count DESC;
>>top10SpamReviewer= LIMIT spamReviewerOrdered 10;
>>SpamTotalData1= JOIN top10SpamReviewer by reviewerID, spamtest by reviewerID;
>>SpamTotalData= FOREACH SpamTotalData1 GENERATE top10SpamReviewer::reviewerID as reviewerID, spamtest::reviewerName as
reviewerName, spamtest::reviewText as reviewText;
```

```
>>STORE SpamTotalData INTO
'gs://dataproc-staging-us-central1-24631865637-6nswal7o/google-cloud-dataproc-metainfo/edab7e1e-e9ba-4097-bace-b003096eb0
fa/cluster-spam-m/SpamDataForTFIDF.json' USING JsonStorage();
```

The next challenge is to concatenate all the rows of different comments and store it in a text file which will be required while performing TFIDF on top 10 spam accounts.
To overcome this, we have taken the aid of python to concatenate the data of each row and create a text file for each reviewer. We have loaded the json file which we had stored in cloud: SpamDataForTFIDF in our python code. Using data model we have joined the reviewText for each reviewer and have written the data in text file.



**Figure 17: Python code to merge all texts for each reviewer**

After this code has generated the 10 text files, I have stored them in the cloud bucket, such that they can be used further by the TF-IDF map reduce code. Now we have the all the data required to perform the further steps.



**Figure 18: Text Files uploaded in Google Cloud Bucket**

## ii. Creating Mappers and Reducers

Python scripts are used in Hadoop to implement TFIDF. There are four mapper program files and three reduction program files in all. There are four stages to the implementation process. Three mappers and three reducers are used in the first phase. In the initial phase we remove the stopwords from the text and clean it further, following that we calculate the word count and the TF. The final phase calculates the TF-IDF and creates a single file containing the spam or ham word lists of 10 users and their TF-IDF scores using the fourth mapper.

One phase's output is used as an input for the following phase. Each text file is provided to the MapReduce.sh script as an input.

This method is applied to each of the ten text files, and the resulting text file for each of the ten text files gets merged at '/home/agnideep_mukherjee2/TFIDF/results/'.

I have also created run shell scripts to integrate the mappers and reducers and provide the input and output path of each phase.

```
C: > Users > 91847 > Downloads > TFIDF MapReduce > $ run-A1KRF81GTI2KKT.sh
 1    hadoop jar /usr/lib/hadoop/hadoop-streaming.jar \
 2    -file /home/agnideep_mukherjee2/TFIDF/mapper_1.py \
 3    -mapper 'python3 mapper_1.py' \
 4    -file /home/agnideep_mukherjee2/TFIDF/reducer_1.py \
 5    -reducer 'python3 reducer_1.py' \
 6    -input gs://dataproc-staging-us-central1-24631865637-6nswal7o/TFIDF_Ham/A1KRF81GTI2KKT.txt \
 7    -output gs://dataproc-staging-us-central1-24631865637-6nswal7o/TFIDF_Ham/out1
 8
 9    hadoop jar /usr/lib/hadoop/hadoop-streaming.jar \
10    -file /home/agnideep_mukherjee2/TFIDF/mapper_2.py \
11    -mapper 'python3 mapper_2.py' \
12    -file /home/agnideep_mukherjee2/TFIDF/reducer_2.py \
13    -reducer 'python3 reducer_2.py' \
14    -input gs://dataproc-staging-us-central1-24631865637-6nswal7o/TFIDF_Ham/out1/ \
15    -output gs://dataproc-staging-us-central1-24631865637-6nswal7o/TFIDF_Ham/out2
16
17    hadoop jar /usr/lib/hadoop/hadoop-streaming.jar \
18    -file /home/agnideep_mukherjee2/TFIDF/mapper_3.py \
19    -mapper 'python3 mapper_3.py' \
20    -file /home/agnideep_mukherjee2/TFIDF/reducer_3.py \
21    -reducer 'python3 reducer_3.py' \
22    -input gs://dataproc-staging-us-central1-24631865637-6nswal7o/TFIDF_Ham/out2/ \
23    -output gs://dataproc-staging-us-central1-24631865637-6nswal7o/TFIDF_Ham/out3
24
25    hadoop jar /usr/lib/hadoop/hadoop-streaming.jar \
```

**Figure 19: Example of a run shell script**

I have upload the 4 mappers and 3 Reducers and Run scripts for the 10 spam text files generated in data processing phase in our Hadoop File system as

```
agnideep_mukherjee2@cluster-spam-m:~$ cd TFIDF
agnideep_mukherjee2@cluster-spam-m:~/TFIDF$ ls
mapper_1.py   reducer_1.py   run-A1KRF81GTI2KKT.sh   run-A2QRXQPHDMFCQV.sh   run-?
mapper_2.py   reducer_2.py   run-A1VAUKYQDX4H2E.sh   run-A34LXQ9YBD2IZT.sh   run-?
mapper_3.py   reducer_3.py   run-A27IN57YAPDJ8S.sh   run-A3AYSYSLHU26U9.sh   sort
mapper_4.py   results        run-A28Q5CQAOGNUM3.sh   run-A3LDPF5FMB782Z.sh
```

**Figure 20: Files uploaded in Hadoop File System**

I have also created a python script 'sort_results.py' for sorting the top 10 results and displaying them in an understandable tabular format.

```
C: > Users > 91847 > Downloads > TFIDF MapReduce > 🐍 sort_results_ham.py
 1    import pandas as pd;
 2    import glob;
 3
 4    folderPath = "/home/agnideep_mukherjee2/TFIDF_Spam/results";
 5
 6    allFiles = glob.glob(folderPath + "/*.txt");
 7
 8    for completeFilePath in allFiles:
 9        filename = completeFilePath.split('/tfidResults/')[1];
10        print('----------------------------------------');
11        print('    '+filename.split('.')[0]+' <- Reviewer ID');
12        df = pd.read_csv(completeFilePath, sep="\t",header=None,names=["word","
13        print('----------------------------------------');
14        df["word"] = df["word"].str.split(" ",n=1,expand=True);
```

**Figure 21: Python code to print all the TF-IDF results in tabular format**

iii. Initiating the process through cloud ssh.

Now I shall run the shell scripts run-<reviewerID> one by one for 10 Spam account reviewers and
Print the data in tabular format using sort_results.py



**Figure 22: TF-IDF results for Spam**

b. TF-IDF of Top 10 ham keywords for top 10 ham accounts

I have followed the same steps and logic that I have previously applied while calculating TF-IDF for top 10
spam keywords in the top 10 spam accounts.

I am attaching the result for the same.

```
----------------------------------
      A64S8V75ITLFG <- Reviewer ID
----------------------------------
                 word   TF-IDF Score
743       badsomeone       0.045709
290      buyingsome        0.017680
142             load       0.014230
448            sorts       0.014230
103    capacitythese       0.009918
300       difference       0.009056
830          backthe       0.008193
43           noticed       0.008193
732           useful       0.007762
880              put       0.007331
```

```
agnideep_mukherjee2@cluster-spam-m:~/TF
------------------------------------
    A2QRXQPHDMFCQV <- Reviewer ID
------------------------------------
             word   TF-IDF Score
56          versa       0.016893
597       between       0.014078
495        period       0.013609
69          canon       0.011732
272      unusable       0.011262
788         handy       0.010793
7          caller       0.010793
566         noise       0.009385
210      changing       0.008916
87     friendliness    0.008447
```

```
------------------------------------
    A3AYSYSLHU26U9 <- Reviewer ID
------------------------------------
             word   TF-IDF Score
400           car       0.021107
340     interfered     0.017684
362       paisley       0.013691
303         until       0.009698
139        rubber       0.008557
712         touch       0.007986
560     available       0.006845
776     recognized     0.006275
623      together       0.005705
845       looking       0.005705
```

```
------------------------------------
    A3LDPF5FMB782Z <- Reviewer ID
------------------------------------
             word   TF-IDF Score
404         phase       0.021192
593          five       0.020400
41          levels       0.014656
1018       amazing       0.013270
1069      moneythe       0.012082
48           menu       0.011487
576      carpetfor       0.009705
1175         major       0.009705
589         email       0.008913
1047     expanding       0.007328
```

```
------------------------------------
    A1KRF81GTI2KKT <- Reviewer ID
------------------------------------
             word   TF-IDF Score
163        moment       0.010471
1155      twisted       0.009043
299        scored       0.006663
6         carrying       0.006663
447       distinct       0.006188
454     fingernail      0.005712
418       voltage       0.005712
1134     processthe     0.005712
663          each       0.005712
730          turn       0.005236
```

```
------------------------------------
    A1VAUKYQDX4H2E <- Reviewer ID
------------------------------------
             word   TF-IDF Score
216       bubbles       0.051813
72            out       0.025907
178     protection     0.020725
50        products       0.018135
167         broke       0.015544
17            cut       0.015544
136          home       0.012953
80           thus       0.012953
123        around       0.012953
97            use       0.012953
```

```
------------------------------------
    A34LXQ9YBD2IZT <- Reviewer ID
------------------------------------
             word   TF-IDF Score
113       bubbles       0.046200
149          once       0.028316
145          full       0.017884
192           one       0.014903
346          very       0.013413
374         under       0.013413
118           few       0.013413
90          looked       0.011923
268        sinkthe       0.010432
256     lightweight    0.010432
```

```
------------------------------------
    A28Q5CQAOGNUM3 <- Reviewer ID
------------------------------------
             word   TF-IDF Score
274        pretty       0.025114
125       thought       0.015982
118          over       0.013699
231           car       0.013699
10         okayish       0.013699
168        phones       0.011416
216          once       0.011416
47          major       0.011416
221     purchased      0.011416
250     purchasing     0.011416
```

```
------------------------------------
    A27IN57YAPDJ8S <- Reviewer ID
------------------------------------
             word   TF-IDF Score
72       peformance      0.037180
79       scratching      0.027885
243          goin       0.024787
398         bunch       0.024012
198        bubble       0.017041
54         donated       0.014717
34        problems       0.013943
177         picky       0.013168
70         mirror       0.010844
60           glue       0.010070
```

```
------------------------------------
    A680RUE1FDO8B <- Reviewer ID
------------------------------------
             word   TF-IDF Score
163        moment       0.010471
1155      twisted       0.009043
299        scored       0.006663
6         carrying       0.006663
447       distinct       0.006188
454     fingernail      0.005712
418       voltage       0.005712
1134     processthe     0.005712
663          each       0.005712
730          turn       0.005236
```

**Figure 23: TF-IDF for Ham Accounts**

# 6) References

1. https://www.learndatasci.com/glossary/tf-idf-term-frequency-inverse-document-frequency/

2. https://pig.apache.org/docs/r0.17.0/api/org/apache/pig/piggybank/storage/CSVExcelStorage.html

3. https://www.geeksforgeeks.org/tf-idf-model-for-page-ranking/

4. https://data.stackexchange.com/stackoverflow/query/new

5. https://stackoverflow.com/questions/20731966/regex-remove-all-special-characters-except-numbers

6. https://www.geeksforgeeks.org/removing-stop-words-nltk-python/