TOPIC :

CUSTOMER COMPLAINTS
ANALYSIS(Banking).


BY :

VISHAL PANDEY

(ROLL NUMBER : 22)


UNDER THE GUIDANCE OF :

PROF. ZAINAB KHAN

# INDEX :

# INTRODUCTION

Consumer-generated content on the Internet continues to grow and impact the banking industry.  The so-called big data analytics approach emphasizes and lever-ages the capacity to collect and analyze data with an unprecedented breadth, depth, and scale to solve real-life problems. In the banking field, there is a growing interest in utilizing user-generated data to gain insights into redressal problems that have not been well understood by conventional methods.

Indeed, big data analytics opens the door to numerous opportunities to develop new knowledge to reshape our understanding of the field and to support decision making in the banking industry. However, while a handful of studies have employed new data sources to tackle important research problems, they were conducted on an ad hoc basis and the application of the big data analytics approach in banking is yet to be well developed and established. The goal of this study is to explore and demonstrate the utility of big data analytics by using it to study core banking management variables and propose a good option for redressal mechanism for customer/consumer complaints and redressals.Previous studies define a complaint as a conflict between a consumer and a business organization in which the fairness of the resolution procedures, the interpersonal communication and behavior, and the outcome of the complaint resolution process are the principal evaluative criteria used by the customer.

In our opinion, a complaint is not necessarily a conflict, however, it can create a conflict between a customer and a business organization, when the answer to the consumer's complaint is not satisfactory.The consumer complaint dataset handles all complaints according to states, products, zip codes, company and sub product categories. Research questions are formulated with the focus on using online customer reviews to enrich our understanding of these constructs. The methodology section details data collection and the text analytical approach utilized to answer the research questions. Findings are then presented and discussed. Finally, the study's contributions to literature and practice as well as directions for future research are discussed.

# IMPLEMENTATION DETAILS

This block of our project basically deals with a core or base software and then the dependencies associated with it so we have divided it in a two block flowchart as follows:
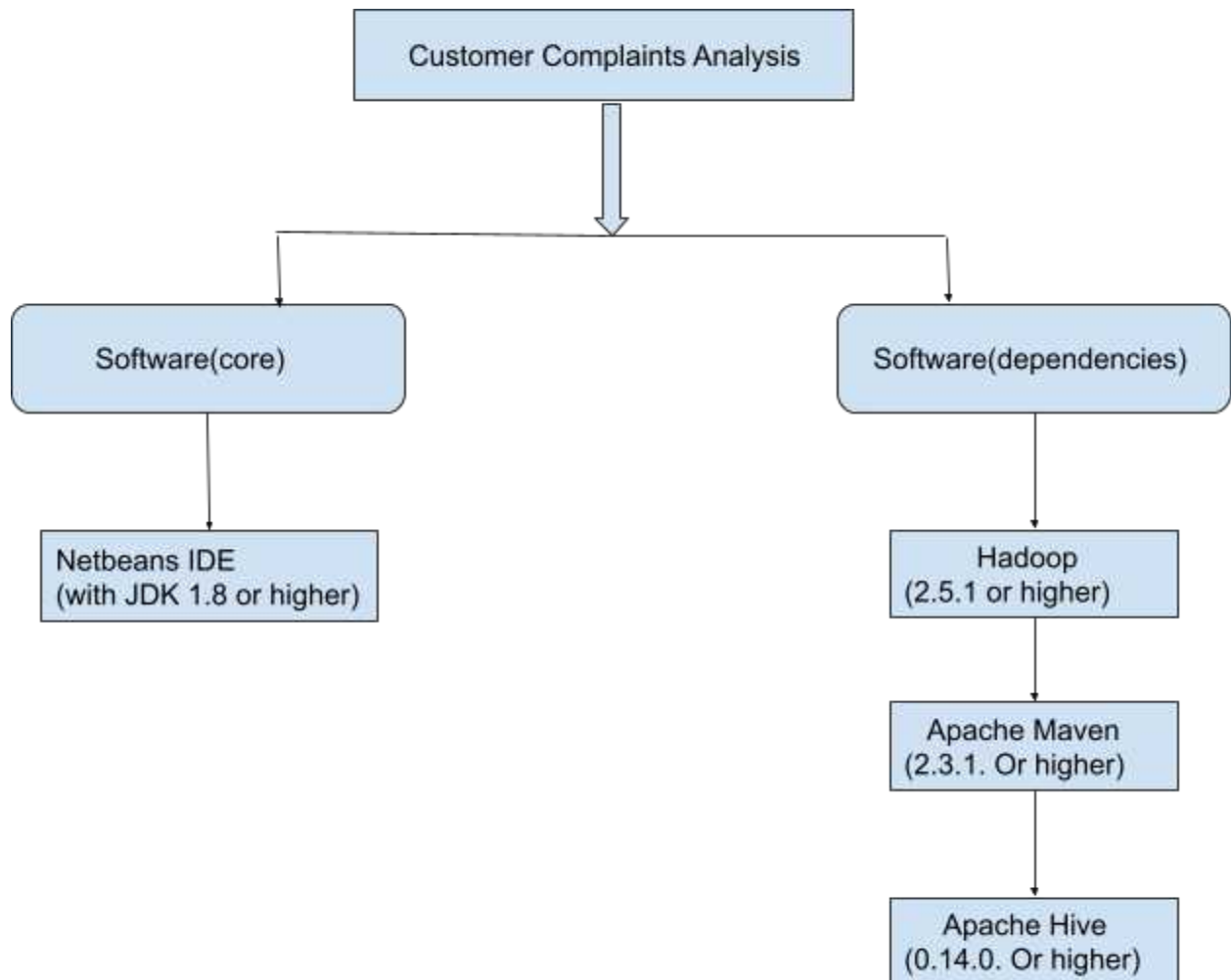


Fig: Flowchart for implementation of CCA

Now further we would like to tell you what exactly needs to be done for smooth running of our project module that is Customer complaints analysis(CCA) in banking industry:

- First and foremost comes a basic hardware requirement to run Netbeans IDE with a version of JDK (java development kit) 1.8 or higher version which requires:

- Any windows version 7 or above
- RAM : 128MB
- Disk Space : 124MB for JRE; 2MB for Java update
- Processor : Minimum pentium 2 266 MHz processor
- Browser : Internet Explorer 9 and above ; Firefox
- NetBeans IDE 8.2 or higher version
- Operating System: Windows 10 / Window 8 / Windows 7 / XP / Vista / Linux
- Memory (RAM): 512MB (for good results)
- Hard Disk: 900MB
- Processor: Intel Pentium IV or above.

- Secondly we have to install all the dependencies mentioned above such as

  - Apache Hadoop (2.5.1 or higher ) and installing it requires basic requirements such as RAM : 1GB, HDD 20GB etc
  - Apache Maven (2.3.1 or higher ) and installing it has prerequisites such as follows JDK : require JDK 1.7 or above to execute - they still allows you to build against 1.3 and other JDK versions ,   Disk : Approximately 10MB is required for the Maven installation itself In addition to that, additional disk space will be used for your local Maven repository. The size of your local repository will vary depending on usage but expect at least 500MB.
  - Apache Hive (0.14.0. Or higher ) also has some prerequisites that need to be taken care of such as Apache Maven .

Next we have to ensure that Hadoop ( 2.5.1 or higher) , and Maven are downloaded and installed as they are mandatory for build. Also the sequence files will be sorted in the by default file system provided by Hadoop which is HDFS.After all this is done we need to run our main project file which has the back end code for running the three modules of our project such as classification , metrics performance and analysis in the last. The back end code is written in java and will run on Netbeans IDE so we need to start it from there after that we have to run some path environment settings to set up java_home , maven home , hadoop home and hdfs so that classification is done properly on the data provided.

Simultaneously there are some .CSV files (Comma-separated values) to test our modules for extensive analysis and classification. These .csv files are taken from the internet with due permission and are used in our module , these are available with us on the webpage link provided. Lastly we need to run maven repositories and their POM (project object model). The Classification takes a bit of time as the hardware we tested it on has some limitations and the speed of execution can be improved with scalability and availability of resources.

Also if the classification and analysis is done once on a particular machine that machine responds faster enough hence improving the performance metrics.The .CSV files would be in a jar format also the links would be provided to pre download it. Next we need to execute some mahout naive bytes classification steps and using mahout and performance matrix analysis too. Also we would require the .CSV files from the CSVReader jar file( the link would be provided as it would act as our dataset). Next we would need Hive to be installed and Hadoop for analysis , The analysis would be shown throughout using geographic maps , graphs,pie charts and bar graphs. These analysis will be with respect to the dataset of the customer complaints and their form of complaints. For example which states have more number of X complaints or which company's complaints are more in Y department of the bank .

# EXPERIMENTAL SETUP AND RESULTS

For the experimental set up the required hardware and software prerequisites have been talked about in the implementation details previously but just to highlight some really important aspects for the projects implementation are JDK installation and any idea that can run java applications or programs for that matter. Also we need Apache maven , Apache hive and Apache mahout for experimenting our project and give viable results and last but not least take care of the environmental variables that you set for different paths in your system as there are some of the most common errors people come across and get frustrated with.

After installation of the required hardware and software prerequisites we need to follow a plan for the experiment which is are divided into two parts and are as follows

- Compilation Instructions : These include a pre defined set path for smoothing data compilation experience and basically it deals with the aspects of getting your repositories cloned , installed and compiled.

- Execution instructions : In this section of experiment we deal with specific queries to get the desired results on the datasets and for that matter CSV (Comma separated values)files which are shot at various terminals such as for hive, hadoop etc.

- [ NOTE: Please download the required CsvReader jar file from the following link
http://javacsv.sourceforge.net/com/csvreader/CsvReader.html ]

# COMPILATION INSTRUCTIONS :

- Project code needs to be cloned from the data repository .Clone the directory and navigate to PROJECT_CODE.
- Install maven if you don't have it.
  - mvn clean install ( Do this after navigating to PROJECT_CODE folder )
- The jar file compiled goes into location ./target/*.jar
  - CONTENTS OF target/*.jar
    - com/
    - com/bigdata/
    - com/bigdata/complaintanalysis/
    - com/bigdata/complaintanalysis/ColumnReducer.class
    - com/bigdata/complaintanalysis/ClassificationAutomator.class
    - com/bigdata/complaintanalysis/StripColumn.class
    - com/bigdata/complaintanalysis/ComplaintsCSVtoSeq.class
    - com/bigdata/complaintanalysis/StatewiseSorter.class

- All required processes should be running for hadoop check using jps command



- Check root directory file content for POM a XML which has all the dependencies that are needed to be compiled

- Classification automator main file which basically calls other java class files inside the source directory



- After going into source java directory we can see a list of java files which are built by maven because they are dependent on each other and required for compilation and execution

- Now do a maven built now by a clean install

- Maven will basically install all the necessities which we require on our system



- We will use hadoop jar command to compile and we will need hadoop packages also and then



This is all in the compilation section of our module, The next section is a continuation of our modules execution after what we have successfully done till now.

# EXECUTION INSTRUCTIONS

In the execution section of the module we are going to deal with the execution of our main classes and running CSV files into our hdfs system for better naive bayes classification and its related confusion matrix in a glance.

- So we can now execute our class file, we can use our main class file as input and also supply an input which is a csv file which is a cleaned dataset.



- This actually creates a lot of log files in the background and the states and banks are split. We also have the sequence files being generated by ComplaintstoCSVsequence generator.

```
StatewiseSorter::sortState() No match exists for state :VT
StatewiseSorter::sortState() No match exists for state :VI
Make new file
StatewiseSorter::sortState() No match exists for state :VI
StatewiseSorter::sortState() No match exists for state :SD
Make new file
StatewiseSorter::sortState() No match exists for state :SD
StatewiseSorter::sortState() No match exists for state :AR
Make new file
StatewiseSorter::sortState() No match exists for state :AR
StatewiseSorter::sortState() No match exists for state :WY
Make new file
StatewiseSorter::sortState() No match exists for state :WY
StatewiseSorter::sortState() No match exists for state :HI
Make new file
StatewiseSorter::sortState() No match exists for state :HI
StatewiseSorter::sortState() No match exists for state :AS
Make new file
StatewiseSorter::sortState() No match exists for state :AS
StatewiseSorter::sortState() No match exists for state :AP
Make new file
StatewiseSorter::sortState() No match exists for state :AP
StatewiseSorter::sortState() No match exists for state :ND
Make new file
StatewiseSorter::sortState() No match exists for state :ND
StatewiseSorter::sortState() No match exists for state :MH
Make new file
StatewiseSorter::sortState() No match exists for state :MH
StatewiseSorter::sortState() No match exists for state :AE
Make new file
StatewiseSorter::sortState() No match exists for state :AE
StatewiseSorter::sortState() No match exists for state :GU
Make new file
StatewiseSorter::sortState() No match exists for state :GU
StatewiseSorter::sortState() No match exists for state :MP
Make new file
StatewiseSorter::sortState() No match exists for state :MP
StatewiseSorter::sortState() No match exists for state :AA
Make new file
StatewiseSorter::sortState() No match exists for state :AA
StatewiseSorter::sortState() No match exists for state :FM
Make new file
StatewiseSorter::sortState() No match exists for state :FM
StatewiseSorter::sortState() No match exists for state :PW
Make new file
StatewiseSorter::sortState() No match exists for state :PW
```

- These Sequence files which are generated by the ComplaintstoCSVsequence generator will then be automatically pushed to our hdfs( hadoop distributed file system).

```
StripColumn::deleteColumn(): Successfully written 2477 rows with TN.csv state
temp_IA.csv
ClassificationAutomator::main(): Created directory data/classification/IA
StripColumn::deleteColumn(): Successfully written 797 rows with IA.csv state
temp_NY.csv
ClassificationAutomator::main(): Created directory data/classification/NY
StripColumn::deleteColumn(): Successfully written 12238 rows with NY.csv state
temp_RI.csv
ClassificationAutomator::main(): Created directory data/classification/RI
StripColumn::deleteColumn(): Successfully written 717 rows with RI.csv state
temp_NE.csv
ClassificationAutomator::main(): Created directory data/classification/NE
StripColumn::deleteColumn(): Successfully written 547 rows with NE.csv state
temp_AZ.csv
ClassificationAutomator::main(): Created directory data/classification/AZ
StripColumn::deleteColumn(): Successfully written 3444 rows with AZ.csv state
temp_PA.csv
ClassificationAutomator::main(): Created directory data/classification/PA
StripColumn::deleteColumn(): Successfully written 6441 rows with PA.csv state
temp_MI.csv
ClassificationAutomator::main(): Created directory data/classification/MI
StripColumn::deleteColumn(): Successfully written 3958 rows with MI.csv state
temp_MT.csv
ClassificationAutomator::main(): Created directory data/classification/MT
StripColumn::deleteColumn(): Successfully written 326 rows with MT.csv state
temp_ND.csv
ClassificationAutomator::main(): Created directory data/classification/ND
StripColumn::deleteColumn(): Successfully written 162 rows with ND.csv state
temp_TX.csv
ClassificationAutomator::main(): Created directory data/classification/TX
StripColumn::deleteColumn(): Successfully written 14090 rows with TX.csv state
temp_WA.csv
ClassificationAutomator::main(): Created directory data/classification/WA
StripColumn::deleteColumn(): Successfully written 3150 rows with WA.csv state
temp_NH.csv
ClassificationAutomator::main(): Created directory data/classification/NH
StripColumn::deleteColumn(): Successfully written 736 rows with NH.csv state
temp_FM.csv
ClassificationAutomator::main(): Created directory data/classification/FM
StripColumn::deleteColumn(): Successfully written 5 rows with FM.csv state
temp_ME.csv
ClassificationAutomator::main(): Created directory data/classification/ME
StripColumn::deleteColumn(): Successfully written 608 rows with ME.csv state
temp_NM.csv
ClassificationAutomator::main(): Created directory data/classification/NM
StripColumn::deleteColumn(): Successfully written 943 rows with NM.csv state
```

- Now hdfs ( hadoop distributed file system) will be called and files are pushed into hadoop, Also these sequence files are hexadecimal encoded.

- After the above operation is done successfully we now have a lot of useful directories inside the hdfs such as states, classification and bank directory. And inside the classification directory we will have a list of states and each state directory will have its chunk-0 file which basically is a sequence file for every state's CSV.

- For example: Let us take New York state (NY) state sequence file which looks hexadecimally encoded and we will use this to generate vectors



- Hexadecimal encoding of the sequence file of the particular state such as New York (NY) as our example is as follows

- Now we need vectors so in Apache mahout now we will generate the vectors, also once the vectorization is generated we will use tfidf vectors to create testing and training datasets



- Basically we have to split the two vectors,datasets and test accordingly so the following command helps us achieving the Split, after the split is done we will do training of the dataset as follows

- This is training the dataset with the use of vectors,training the module; We basically have to first train the model and secondly test it, So the following one is the training dataset



- The next step is to run the vectors and test the dataset which we have after running the training dataset, So now we will be testing the dataset as follows

- After training now we will generate confusion matrix by doing a naive bayes classification on the train set as well as the test dataset which goes as follows

- Using Apache Mahout we will do the classification which is also called as the Naive Bayes classification , In Naive bayes classification we will train the training dataset, which usually gives good results and also provides a higher accuracy of around 95 to 100% but once we run it on the testing dataset it gives clearly low accuracy between 85 to 90% which is expected and measured.



- Now we will get to the interesting part of our module which is the Confusion matrix , Confusion matrix actually has classifiers such as credit card, debit card, mortgage, loan, etc which are basically some of the functions of the banking industry and has a lot of complaints regarding it. So basically it simplifies and classifies and we can see those trends by the classifiers it has  for ex: complaints related to credit card, loan,etc with respect to our project and input.

```
14/12/22 00:27:26 INFO mapred.JobClient:    FileSystemCounters
14/12/22 00:27:26 INFO mapred.JobClient:        FILE_BYTES_READ=25928128
14/12/22 00:27:26 INFO mapred.JobClient:        FILE_BYTES_WRITTEN=26484859
14/12/22 00:27:26 INFO mapred.JobClient:    Map-Reduce Framework
14/12/22 00:27:26 INFO mapred.JobClient:        Map input records=10353
14/12/22 00:27:26 INFO mapred.JobClient:        Physical memory (bytes) snapshot=0
14/12/22 00:27:26 INFO mapred.JobClient:        Spilled Records=0
14/12/22 00:27:26 INFO mapred.JobClient:        Total committed heap usage (bytes)=77070336
14/12/22 00:27:26 INFO mapred.JobClient:        CPU time spent (ms)=0
14/12/22 00:27:26 INFO mapred.JobClient:        Virtual memory (bytes) snapshot=0
14/12/22 00:27:26 INFO mapred.JobClient:        SPLIT_RAW_BYTES=161
14/12/22 00:27:26 INFO mapred.JobClient:        Map output records=10353
14/12/22 00:27:28 INFO test.TestNaiveBayesDriver: Complementary Results:
=====================================================
Summary
-----------------------------------------------------
Correctly Classified Instances       :       9895     95.5762%
Incorrectly Classified Instances     :       458      4.4238%
Total Classified Instances           :       10353
=====================================================
Confusion Matrix
-----------------------------------------------------
a       b       c       d       e       f       g       h       i       <--Classified as
965     0       0       0       0       0       0       0       0       |  965     a  =  Bank account or service
0       265     0       0       0       0       0       0       57      |  322     b  =  Consumer loan
0       0       2451    154     0       0       242     0       0       |  2847    c  =  Credit card
0       0       0       2682    0       0       0       0       0       |  2682    d  =  Credit reporting
0       0       0       0       1853    0       0       0       0       |  1853    e  =  Debt collection
0       0       0       0       0       136     0       0       0       |  136     f  =  Money transfers
0       0       0       0       0       0       612     0       0       |  612     g  =  Mortgage
0       0       0       0       0       0       0       43      5       |  48      h  =  Payday loan
0       0       0       0       0       0       0       0       888     |  888     i  =  Student loan
=====================================================
Statistics
-----------------------------------------------------
Kappa                                  0.9439
Accuracy                               95.5762%
Reliability                            85.7972%
Reliability (standard deviation)       0.3091
Weighted precision                     0.9636
Weighted recall                        0.9558
Weighted F1 score                      0.9563

14/12/22 00:27:28 INFO driver.MahoutDriver: Program took 7145 ms (Minutes: 0.11908333333333333)
blitzavi89@blitzavi89-Lenovo-Ideapad-Flex-14:~/BigData_Project/big_data_analytics/PROJECT_CODE$
```

- Now we will repeat the analysis on the testing dataset to validate the results for our module and we can see the confusion matrix for our classification performed on the testing dataset.

```
0       265     0       0       0       0       0       0       57      |  322     b  =  Consumer loan
0       0       2451    154     0       0       242     0       0       |  2847    c  =  Credit card
0       0       0       2682    0       0       0       0       0       |  2682    d  =  Credit reporting
0       0       0       0       1853    0       0       0       0       |  1853    e  =  Debt collection
0       0       0       0       0       136     0       0       0       |  136     f  =  Money transfers
0       0       0       0       0       0       612     0       0       |  612     g  =  Mortgage
0       0       0       0       0       0       0       43      5       |  48      h  =  Payday loan
0       0       0       0       0       0       0       0       888     |  888     i  =  Student loan
=====================================================
Statistics
-----------------------------------------------------
Kappa                                  0.9439
Accuracy                               95.5762%
Reliability                            85.7972%
Reliability (standard deviation)       0.3091
Weighted precision                     0.9636
Weighted recall                        0.9558
Weighted F1 score                      0.9563

14/12/22 00:27:28 INFO driver.MahoutDriver: Program took 7145 ms (Minutes: 0.11908333333333333)
blitzavi89@blitzavi89-Lenovo-Ideapad-Flex-14:~/BigData_Project/big_data_analytics/PROJECT_CODE$ ~/mahout-trunk/bin/mahout testnb -i test-vectors -m model -l labelindex -ow -o NY-testin
g -c
MAHOUT_LOCAL is set, so we don't add HADOOP_CONF_DIR to classpath.
MAHOUT_LOCAL is set, running locally
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/blitzavi89/mahout-trunk/examples/target/mahout-examples-1.0-SNAPSHOT-job.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/blitzavi89/mahout-trunk/examples/target/dependency/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]
14/12/22 00:27:38 WARN driver.MahoutDriver: No testnb.props found on classpath, will use command-line arguments only
14/12/22 00:27:39 INFO common.AbstractJob: Command line arguments: {--endPhase=[2147483647], --input=[test-vectors], --labelIndex=[labelindex], --model=[model], --output=[NY-testing],
--overwrite=null, --startPhase=[0], --tempDir=[temp], --testComplementary=null}
14/12/22 00:27:39 INFO common.HadoopUtil: Deleting NY-testing
14/12/22 00:27:39 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
14/12/22 00:27:40 INFO input.FileInputFormat: Total input paths to process : 1
14/12/22 00:27:40 INFO filecache.TrackerDistributedCacheManager: Creating filemodel in /tmp/hadoop-blitzavi89/mapred/local/archive/4328442496935006430_104069929_1886782980-work-1831653
693252742768 with rwxr-xr-x
14/12/22 00:27:40 INFO filecache.TrackerDistributedCacheManager: Cached model as /tmp/hadoop-blitzavi89/mapred/local/archive/4328442496935006430_104069929_1886782980/filemodel
14/12/22 00:27:40 INFO filecache.TrackerDistributedCacheManager: Cached model as /tmp/hadoop-blitzavi89/mapred/local/archive/4328442496935006430_104069929_1886782980/filemodel
14/12/22 00:27:41 INFO mapred.JobClient: Running job: job_local1754081869_0001
14/12/22 00:27:41 INFO mapred.LocalJobRunner: Waiting for map tasks
14/12/22 00:27:41 INFO mapred.LocalJobRunner: Starting task: attempt_local1754081869_0001_m_000000_0
14/12/22 00:27:41 INFO mapred.util.ProcessTree: setsid exited with exit code 0
14/12/22 00:27:41 INFO mapred.Task:  Using ResourceCalculatorPlugin : org.apache.hadoop.util.LinuxResourceCalculatorPlugin@7d456c3b
14/12/22 00:27:41 INFO mapred.MapTask: Processing split: file:/home/blitzavi89/BigData_Project/big_data_analytics/PROJECT_CODE/test-vectors/part-r-00000:0+126787
```

- Similarly we can do the classification for banks which measures priorities for different products such as credit card, loans, etc. So banks can use these results in near future in order to calculate the priority and importance for tackling specific issues raised by customers. The classification for banks is shown below for example



- Apache hive was used in order to extract meaningful information from a given dataset in hive, But first we need to create an empty table Cdata for complaint data and fire describe tablename query which gives schema of the table as below in hive

- Now lets insert data into the table, as data is stored locally in a CSV file we will simply instruct hive to copy that local file into the hive table



- Now the table is ready to accept queries so query the table with a query which asks top 5 companies who have highest percentage of closed complaints by customers and have at least 10 closed complaints we have also exported the data in a csv file for future purposes

- The working of the above query goes as the sub query A gives the company name and total number of complaints closed by the company after that we join the table with the result of sub query B which gives company name and number of closed complaints disputed by the customer dividing the result of second with the result of the first we get the percentage of the closed complaints that were disputed ,Results here gives out top 5 companies which have maximum percentage of customers dis-satisfied by their solutions. These results thus are very important from a company perspective as well because it helps them to judge their performance.



- Some additional analysis done via tableau software which gives insights in a more representative manner

- Product-wise analysis

Product



Product and count of Product. Color shows details about Product. Size shows count of Product. The marks are labeled by Product and count of Product.

- Mortgage-issue analysis

Issue - Mortgage



Count of Issue for each Issue. Color shows details about Issue. The data is filtered on Product, which keeps Mortgage.

- Company wise analysis



| Bank of America 57,992 | Experian 33,912 | Capital One 16,439 | Nationstar Mortgage 13,940 | | U.S. Bancorp 10,086 | Ditech |
|---|---|---|---|---|---|---|
| | | PNC Bank N.A. 7,306 | | Fifth | | Ally |
| Wells Fargo & Company 43,956 | TransUnion Intermediate Holdings, Inc. 28,248 | Encore Capital Group | | | | |
| | | HSBC North America | | | | |
| JPMorgan Chase & Co. 35,650 | Citibank 27,535 | Amex 6,193 | | | | |
| | | SunTrust Banks, Inc. | | | | |
| | | Discover 5,076 | | | | |
| Equifax 35,439 | Ocwen 21,780 | TD Bank US Holding | | | | |
| | | Select Portfolio | | | | |

- Year wise analysis



year wise

Date received

The trend of count of Date received for Date received Year.

- Timely response

**Timely Response**



Count of Timely response? for each Timely response?. Color shows details about Timely response?.

- Maximum issues registered state wise

- Complaints submitted via

submitted via



Count of Submitted via for each Submitted via. Color shows details about Submitted via. Details are shown for Submitted via.

# ANALYSIS OF THE RESULTS

- Analysis of the results we had are extremely encouraging as we took the datasets of the Consumer Financial Protection Bureau (CFPB) of the states and made it valuable by classifying it with regards to specific inputs and outputs which can help the current scenario as according to us there are more frustrated people with the current redressal mechanism adopted by the banking industry .

- So talking about the analysis of the results we had with respect to classification done on Apache Mahout with Naive bayes classification we found out that the module had a higher accuracy of about 95 to 100% with the training dataset and the accuracy of our module classification with the testing dataset had an accuracy of 85 to 90% which is shown below and was predicted which in turn tells us about the analysis was good.

```
14/12/22 00:27:43 INFO mapred.JobClient:    FileSystemCounters
14/12/22 00:27:43 INFO mapred.JobClient:        FILE_BYTES_READ=25331649
14/12/22 00:27:43 INFO mapred.JobClient:        FILE_BYTES_WRITTEN=25641153
14/12/22 00:27:43 INFO mapred.JobClient:    Map-Reduce Framework
14/12/22 00:27:43 INFO mapred.JobClient:        Map input records=1885
14/12/22 00:27:43 INFO mapred.JobClient:        Physical memory (bytes) snapshot=0
14/12/22 00:27:43 INFO mapred.JobClient:        Spilled Records=0
14/12/22 00:27:43 INFO mapred.JobClient:        Total committed heap usage (bytes)=60817408
14/12/22 00:27:43 INFO mapred.JobClient:        CPU time spent (ms)=0
14/12/22 00:27:43 INFO mapred.JobClient:        Virtual memory (bytes) snapshot=0
14/12/22 00:27:43 INFO mapred.JobClient:        SPLIT_RAW_BYTES=160
14/12/22 00:27:43 INFO mapred.JobClient:        Map output records=1885
14/12/22 00:27:43 INFO test.TestNaiveBayesDriver: Complementary Results:
=======================================================
Summary
-------------------------------------------------------
Correctly Classified Instances          :       1679      89.0716%
Incorrectly Classified Instances        :        206      10.9284%
Total Classified Instances              :       1885
=======================================================
Confusion Matrix
-------------------------------------------------------
a    b    c    d    e    f    g    h    i    <--Classified as
633  0    0    0    0    0    0    0    0    |  633    a  = Bank account or service
0    175  0    0    0    0    0    0    57   |  232    b  = Consumer loan
0    0    871  84   0    0    65   0    0    |  1020   c  = Credit card
0    0    0    0    0    0    0    0    0    |  0      d  = Credit reporting
0    0    0    0    0    0    0    0    0    |  0      e  = Debt collection
0    0    0    0    0    0    0    0    0    |  0      f  = Money transfers
0    0    0    0    0    0    0    0    0    |  0      g  = Mortgage
0    0    0    0    0    0    0    0    0    |  0      h  = Payday loan
0    0    0    0    0    0    0    0    0    |  0      i  = Student loan
=======================================================
Statistics
-------------------------------------------------------
Kappa                               0.8215
Accuracy                            89.0716%
Reliability                         26.0823%
Reliability (standard deviation)    0.424
Weighted precision                  1
Weighted recall                     0.8907
Weighted F1 score                   0.9401

14/12/22 00:27:43 INFO driver.MahoutDriver: Program took 4859 ms (Minutes: 0.08098333333333334)
blitzavi89@blitzavi89-Lenovo-Ideapad-Flex-14:~/BigData_Project/big_data_analytics/PROJECT_CODES
```

- Performance metric analysis calculation for various categories of issues that are actually entered by the customer are parameters that help us better the calculator and up the game of redressal mechanism. When a customer is entering a specific complaint with the database, he can have issues with credit card, debit card, etc so based on the priority values that are being calculated by the code we can figure out which issue is supposed to be given a higher priority while a customer actually enters his specific complaint allowing him to a better redressal experience. Here in this performance analysis everything ran as predicted keeping the resources in mind.

```
25      static double exponent = 2.718282;
26
27      static HashMap<String, Integer> productList;
28      static HashMap<String, Integer> issueList;
29
30      static HashMap<String, Integer> productCompanyResponseCount; /*Hashes a concat of 2 strings - product and response*/
31      static HashMap<String, Integer> productTimelyResponseCount;
32      static HashMap<String, Integer> productConsumerDisputedCount;
33
34      static HashMap<String, Integer> issueCompanyResponseCount; /*Hashes a concat of 2 strings - issue and response*/
35      static HashMap<String, Integer> issueTimelyResponseCount;
36      static HashMap<String, Integer> issueConsumerDisputedCount;
37
38
39      private void initLists() {
40
41          productList = new HashMap<String, Integer>();
42          issueList = new HashMap<String, Integer>();
43          productCompanyResponseCount = new HashMap<String, Integer> ();
44          productTimelyResponseCount = new HashMap<String, Integer>();
45          productConsumerDisputedCount = new HashMap<String, Integer>();
46
47          issueCompanyResponseCount = new HashMap<String, Integer> ();
48          issueTimelyResponseCount = new HashMap<String, Integer> ();
```

Markers  Properties  Servers  Data Source Explorer  Snippets  Console

ProblemClustering [Java Application] /usr/lib/jvm/java-7-openjdk-amd64/bin/java (Dec 21, 2014, 10:01:55 PM)
Average scores based on products
Mortgage >> 126.91634337466999
Money transfers >> 145.74958096900264
Consumer loan >> 143.51746888360364
Student loan >> 140.04344170275303
Debt collection >> 138.42967745770443
Payday loan >> 144.81518983014308
Credit card >> 120.58988461872354
Bank account or service >> 127.48735903172147

**Here are the metrics for various categories of issues / products**

- The performance metric analysis calculator which did calculation for the current module can of course be further improvised by using map reduce version of it. It will also help in reducing response time of the query, better availability and seamless experience with the whole complaint to resolution cycle reduced to its limit.

```
25    static double exponent = 2.718282;
26
27    static HashMap<String, Integer> productList;
28    static HashMap<String, Integer> issueList;
29
30    static HashMap<String, Integer> productCompanyResponseCount; /*Hashes a concat of 2 strings - product and response*/
31    static HashMap<String, Integer> productTimelyResponseCount;
32    static HashMap<String, Integer> productConsumerDisputedCount;
33
34    static HashMap<String, Integer> issueCompanyResponseCount; /*Hashes a concat of 2 strings - issue and response*/
35    static HashMap<String, Integer> issueTimelyResponseCount;
36    static HashMap<String, Integer> issueConsumerDisputedCount;
37
38
39    private void initLists() {
40
41        productList = new HashMap<String, Integer>();
42        issueList = new HashMap<String, Integer>();
43        productCompanyResponseCount = new HashMap<String, Integer> ();
44        productTimelyResponseCount = new HashMap<String, Integer>();
45        productConsumerDisputedCount = new HashMap<String, Integer>();
46
47        issueCompanyResponseCount = new HashMap<String, Integer> ();
48        issueTimelyResponseCount = new HashMap<String, Integer> ();
```

Markers  Properties  Servers  Data Source Explorer  Snippets  Console ⊠

ProblemClustering [Java Application] /usr/lib/jvm/java-7-openjdk-amd64/bin/java (Dec 21, 2014, 10:01:55 PM)
```
Average scores based on products
Mortgage >> 126.91634337466999
Money transfers >> 145.74958096900264
Consumer loan >> 143.51746888360364
Student loan >> 140.04344170275303
Debt collection >> 138.42967745770443
Payday loan >> 144.81518983014308
```

**Calculating the average CPM of various financial products for a subset of our data. Performance can be improved with the MapReduce version of the same algorithm, which we plan to release in the future**

# CONCLUSION

- From the above mentioned analysis, it is seen that the credit card complaints are increasing since the last 5 years. From the business point of view, the home loan objections can be decreased to a degree if the banks' records are frequently checked by the powers and the banks ought to be advised to deliver and submit reports of their exercises all the time barring the clients' private and classified information.

- Steps should be taken to verify all documents and avoid customers dissatisfaction. During this project we learned and improved our analytical skills. Working on Hadoop Distributed File System helped us gain knowledge on Hive Query Language.

- Meanwhile we were able to represent queries into graphical form. In our detailed analysis, we were able to analyze and represent the maximum number of complaints in California with the majority of complaints being of credit card Complaints. Also, through our analysis, we found which banks are complained about the most and via which web are the complaints recorded. Also it came to our notice that all complaints were timely responded.

- If talking about very specific analysis trend and details drawn from it we can also conclude the following
    - Mortgage product and other mortgage sub-product complaints are maximum in number.
    - California has the highest number of complaints.
    - Most of the complaints are registered with Bank of america.
    - Consumer complaints are increasing year by year.

- This brought us to conclusion providing useful aspects in regards to recurrence of complaints, complaints from different places, about diverse products, sub-products and companies which make us aware and take necessary steps to avoid them.

- With our new metric system, banks can relatively prioritize the complaints to resolve as we proposed a system to evaluate a good redressal mechanism for banking customers/consumers complaints based on the analysis.

- By analysing and processing the text, the system can make the reviews into a graph that is easy to read and understand. Performed data analysis on the data sets, to give a detailed overview of the banks performance from a customer sentiment perspective.

- Developed a novel metric system that assigns priorities to customers complaints. This helps banks prioritize customers problems on specific constraints such as response time, etc.

# FUTURE ENHANCEMENT

- Resolution Methodology Recommender.Build a recommender engine that can derive the best "first response" for a complaint.

- More rigorous data analysis and research into complaint resolution methodologies required.

- Quantitatively gauge impact of various classes of complaints on various products to gain insights into customer outlook towards a specific product.

- Seek to make a better performance metric analysis calculator with the map reduce version of what was used in this module of customer complaints analysis with respect to banking

# PROGRAM CODE

● Program code of this project has some major files coded in java and xml both as all the main class files are being called by a file called ClassificationAutomator which basically calls all the useful and dependent java class files and the main class too.

- Then comes one of the main java files named ComplaintsCSVtoSeq java file it is the heart of our project as it generates sequence files and pushes it automatically to the hdfs.





- A Project Object Model or POM is the fundamental unit of work in Maven. It is an XML file that contains information about the project and configuration details used by Maven to build the project. It contains default values for most projects. Examples for this are the build directory, which is target; the source directory, which is src/main/java;

the test source directory, which is src/test/java; and so on. When
executing a task or goal, Maven looks for the POM in the current
directory. It reads the POM, gets the needed configuration information,
then executes the goal.

- Then come the secondary files mandatory for smooth execution such as vectors and bank and states sorter files also we have not uploaded all the files for all files we have drop our github repository link where you can find the whole project code and some more valuable information.

Following is the code for Bankwise Sorting

Next is the State wise Sorter file





- The link for the whole project code and some more files related to analysis and execution can be found at
  [ https://github.com/Vshal-P/CCA ]