

# FINAL REPORT

## Sentiment Analysis

KUNWAR RANANJAY SINGH

2017B4A70504G

f20170504@goa.bits-pilani.ac.in



**BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE, PILANI  
GOA CAMPUS**

## CONTENTS

Problem	3
Related Work	4
Approach and Solution	5
Code and Output	12
Data Collection	39
Evaluation of Solution	40
Conclusion	51
Future Work	52
References	53

## **Problem**

Aim is to classify social media content using data science into positive, negative and neutral categories. It's also known as sentiment analysis.

For example : “enjoyed travelling today. Will visit again” can be classified as positive whereas “lot of rain today. Depressing !!!” can be classified as negative.

Social media has a lot of impact on the human brain and using technology which can classify posts into various categories can be used to filter out those which have a negative impact. It can be also used by companies to only present posts of a certain category that the user has selected(for example filtering out negative posts). Customer posts on websites can be used to improve products. It can be used by governments to reduce spread of misinformation and can also be used for fake news detection [1].

## Related Work

Let's see some of the important points from the work done by other researchers.

Majority of the work in the field of sentiment analysis has been done on twitter data. In previous works it has been observed that an algorithm provides more accurate results on short tweets. In one of the studies done on COVID-19 tweets, Naive Bayes is better compared to Logistic Regression in terms of accuracy for both short and long tweets. Each model has its own strengths and weaknesses. Linear Regression is much easier to understand conceptually however it's prediction accuracy can be affected by outliers. Logistic Regression can be used to avoid overfitting compared to Polynomial Regression however it depends on assumption that data is linear. For example in tweets like "I am feeling sad today", "happy to visit Goa", we can see words like "sad", "happy" which would be present in a lot of negatively and positively classified tweets respectively and thus it becomes easier to classify them using logistic regression. However in tweets like "Movie was bad even though all the cast and crew members were great", we can see it has words like "bad", "great" which are the only words impacting sentiment of this tweet. Overall it has negative sentiment however by just going with the meaning of words it would be neutral. The Naive Bayes classifier has an assumption that features are independent. K-Nearest neighbour can be used to classify data into multiple categories like fear, sadness, anger, disgust however it's computationally inefficient for large data sets. According to some of the word clouds (where the size of a word is proportional to its frequency) it was observed that "trump", "virus", "stock" were present in large frequency in a lot of covid-19 tweets. This observation can be utilised to make more sophisticated preprocessing techniques. For example if the majority of tweets containing word "stock" have negative sentiment then while classifying we can have additional negative weightage to tweets having word "stock" provided we are classifying tweets during COVID-19 period [2].

Some of the common pre-processing techniques used are removing URLs, hash tags, usernames, correction of spelling, removing emoticons with their sentiment, removing punctuation, stop words, removing non-english words in tweets. Most work didn't use unsupervised machine learning techniques because we need some kind of mapping of tweets to some quantity which can be utilised for sentiment analysis. Just putting an English word won't give any meaning to the unsupervised algorithms. Lexicon based approaches in which semantic orientation of words and phrases are used seem to have less accuracy compared to machine learning techniques. As expected, accuracy of any ML algorithm is better when the training set is large. It has been observed that bigram which refers to a combination of two words can increase the accuracy since putting bigrams like "not sad", "not good" (they correspond to positive and negative sentiment, respectively) in calculation would reduce the effect of unigrams "sad", "good" (they correspond to negative and positive sentiment, respectively). In one of the experiments accuracy increased to 76.44 from 74.56 percent by using bigrams [3].

Lexicon based approaches can be used for unsupervised learning, but it has its own challenges. Polarity of the words in the sentence can be used to determine the sentiment but natural languages are complex and it won't consider things like negation.

Taking account of slang and emoticons in the sentences can improve the sentiment classification accuracy of the tweets. One study has shown that the presence of emoticons in tweets increased the accuracy from 79% to 85%[4].

## **Approach and Solution**

Focus will be on tweets since they have a character limit of 280 [5]. This is because as the number of characters increases the post becomes more vague and less objective in terms of classifying into positive or negative whereas in tweets it's mostly inclined toward some emotion due to character limit.

I started with the most naive solution. Following steps were followed while coding the naive solution:

1. Splitting the tweets in tokens using space as a delimiter.
2. Converting the words into lower case.
3. Ignoring the words containing characters other than alpha-numeric.
4. Using a dictionary to map a token to count of sentiment (-1 is added if the token belongs to a tweet which has negative sentiment else +1 is added) . In other words we are creating a frequency table for each token where the value added while iterating the tweets is negative if the tweet belongs to negative sentiment else positive.
5. Splitting the dataset into 80-20 for training and testing purposes respectively. Two vectors are created, the first one contains the sum of the dictionary value of tokens in the tweet and the second contains corresponding sentiment (0 for negative, 1 for positive).
6. Logistic Regression and Support Vector Machine are used individually for training and then predicting the sentiment for test data.
7. Confusion matrix and accuracy is printed.

After that I thought of ways (on paper) to clean and preprocess the data in a better way, better and different ways of representing the data mathematically and using more algorithms and techniques. I implemented them in python.

**There were four major steps :**

1. Choosing the datasets.
2. Preprocessing and cleaning of the data.
3. Representing data mathematically.
4. Using algorithms for training and testing.

### **Choosing the datasets**

Two different datasets were chosen to test the performance. First dataset contains random tweets (they don't belong to some specific domain). From now on we will refer to the first dataset as dataset-1. Second dataset contains tweets related to US airlines and we will refer to it as dataset-2.

### **Preprocessing and cleaning of the data**

Following steps represents the whole process in preprocessing and cleaning of the data :

1. Dataframe rows are shuffled. It's possible that a lot of negative or positive tweets are concentrated at some consecutive set of rows in the dataset which can affect training process and test results. To avoid such a situation rows were shuffled.
2. Columns other than those having tweets or sentiment are removed. For example, a column having the date of tweet has no role in determining the sentiment.
3. Sentiment column values are mapped to +1 if tweet was having positive sentiment else -1 if tweet was having negative sentiment. Algorithms work with numbers and thus it's required to map sentiment to +1 or -1.
4. Emojis are replaced with their meaning. For example, happy emoji is replaced with the word "happy".
5. Contractions are expanded. For example "won't" is converted to "will not". This is helpful because there are sentences with the word "will not" and it would be helpful to consider their meaning similar to "won't".
6. Hashtag from words is removed. For example "#summer" is replaced with "summer". Words attached with hash convey meaning and thus it's important to not remove them completely.
7. "RT" representing retweet is removed. "RT" doesn't have any effect on sentiment and thus it should be removed.
8. Words in tweets having double forward slashes "//" are removed. Usually words in tweet with '#' or having underscores convey some meaning. However words with "//" are rare and they usually represent some website which has no role in determining the sentiment.
9. Positions of words having '!' (exclamation mark) are noted and '!' are removed. Finally Words having characters other than lower case, upper case alphabet, underscore are removed. Words beginning with '@' for example '@daniel' represent some other profile and thus it is not useful for determining the sentiment. Exclamation marks usually represent the extreme end of emotion and thus keeping track of words attached with an exclamation mark at the end can be useful for final sentiment calculation.
10. Tweets are changed to lowercase. "Not" and "not" convey the same meaning and thus we need to change all words into either lowercase or uppercase.
11. Words having underscores are splitted. For example "good\_day\_today" is changed to "good day today". Words having underscore are found frequently in tweets (mostly with hash at the front like #good\_day\_today).
12. Spell corrections are done. For example "wowwwwww" can be changed to "woww" (i.e removing multiple repeating characters) using regex and then further corrected to "wow" with spell correction packages.
13. Stop words other than "not" are removed. "is", "a" are some examples of stop words. They usually don't convey any meaning from sentiment analysis perspective. However if they are not removed from data then they can skew the results. For example if "a" is present in a lot of negative tweets in training data and in testing data if it's present in some tweet having positive sentiment, then it will try to push the result toward negative sentiment. "not" is not removed because it acts as negation and useful in cases like "not happy" where it changes the overall meaning.
14. Lemmatization of tweets is done. Stemming and Lemmatization can be used to shorten the words and group together similar words. For example in "Saurabh is Intern", "Saurabh is doing internship", both "intern" and "internship" convey the same thing from sentiment

analysis perspective and thus changing “internship” to “intern” will give better results. Lemmatization was used because it takes care of “part of speech tagging” i.e it considers the context and then makes the word short accordingly. For example it would be bad if both “international” and “internship” is reduced to “intern”. If we take another example to understand the importance of part of speech tagging : in “The man fans the flame”, “fans” is a verb whereas in “The fans watch the race” “fans” is a noun.

15. Those words which are having ‘!’ at the end of them, they are replaced with twice repetition of themselves. For example ‘happy!!’ replaced with ‘happy happy’.
16. Tuples are appended at the end of the tweets. For example, for the tweet : “not good today” and n=2, the tweet will be changed to “not good today \_not\_good\_good\_today”. There are words which often act as negation and thus it becomes important to consider the meaning of consecutive words as a group. For example in “not happy today”, “not” is acting as negation and we don’t need to see “not” and “happy” separately. If we have made tuples, then “\_not\_happy” will be there in the final modified tweet and thus it might prevent from displaying false sentiment. If we take another example : “This comedy movie wasn’t that great. It didn’t make the audience happy”. Here we can see that it’s a negative sentiment. However it has words like “great”, “happy” which would convey positive sentiment. Words that decided the sentiment are “wasn’t” and “didn’t”. They act as negation. However “wasn’t”, “didn’t” are two word distance apart from “great”, “happy” and bigrams won’t help here. We need to consider 3 to 4 words together in a lot of cases to get true sentiment and thus I have considered upto 4-grams(tuples of size upto 4).
17. Dataframe having sentiment column and tweet column is saved as csv file which can be used further for testing with different models. All the previous process takes a lot of time (around 4 hours on my local computer for 20000 tweets and tuples size 4) and thus it would save time if we keep the modified data frame saved so that we don’t need to create it every time we use it in our model.

### Representing data mathematically

Following three different techniques were used to represent the data mathematically so that it can be used by algorithms :

1. Positive-Negative frequency. We create a dictionary where key is a word and value is a pair. First value in the pair represents the number of tweets in training datasets which are positive and that word is present in it. Second value in the pair represents the number of tweets in training datasets which are negative and that word is present in it. In the final dataframe there are three columns. First column represents the sum of positive values of all words in the corresponding tweet at that row index. Second column represents the sum of negative values of all words in the corresponding tweet at that row index. Third column is the corresponding sentiment. Let’s understand it with example of three tweets where we choose 2 tweets for training and 1 tweet for testing :

tweet	sentiment
-------	-----------

"i am happy today"	1
"i am sad today"	-1
"i am feeling happy"	1

Following dictionary would be created from first two rows :

word	positive sum	negative sum
i	1	-1
am	1	-1
happy	1	0
today	1	-1
sad	0	-1

Final dataframe where tweet column is replaced by two columns :

positive sum	negative sum	sentiment
4	-3	1
3	-4	-1
3	-2	1

Now we can use the first two rows for training and the final row for testing.

2. Bag of Words. Here columns are words (number of columns is equal to number of unique words). In each row and for each word we write the count of the number of times that word occurs in the corresponding tweet at the same row. An example will make this clear :

tweet	sentiment
"i am happy today"	1
"i am sad today"	-1
"i am feeling happy"	1

Corresponding data frame representing bag of words and last column as sentiment :

i	am	happy	today	sad	feeling	sentiment
1	1	1	1	0	0	1
1	1	0	1	1	0	-1
1	1	1	0	0	1	1

3. TF-IDF. Term Frequency (TF) of a word in a sentence is defined as  $\frac{c}{L}$  where  $c$  is frequency of that word in the sentence and  $L$  is the total number of words in the sentence. For example in sentence “i like to travel and i like to listen music”, TF of word “like” is  $\frac{2}{10}$ . Inverse Document Frequency (IDF) is defined as  $\log(\frac{N}{n})$  where  $N$  is total number of sentences and  $n$  is the number of sentences in which that word occurs. For example if we are given two sentences : “i like to travel and i like to listen music”, “i am sad because we lost”. Then the IDF of “like” is  $\log(\frac{2}{1})$  since it occurs only in the first sentence and there are a total two sentences. TF-IDF of a word = TF\*IDF. For example in previous two sentences, TF-IDF of “like” in first sentence is  $\frac{2}{10} * \log(\frac{2}{1}) = 0.0602$ . TF-IDF of “i” in first sentence is  $\frac{2}{10} * \log(\frac{2}{2}) = 0$ . High TF-IDF of a word in a sentence implies that it occurs in that sentence more times compared to other sentences. Word “like” occurs only in the first sentence whereas word “i” occurred in both of them and thus word “like” has more TF-IDF value in the first sentence compared to “i”. Our intuition also says to give “like” more importance compared to “i” while finding the sentiment.

In our dataframe we treat the tweets as sentences and columns are words. Then in each row we write TF-IDF value for each word for the corresponding tweet.

Let's understand it with help of an example :

tweet	sentiment
“i am happy today”	1
“i am sad today”	-1
”i am feeling happy”	1

The corresponding data frame having TF values is :

i	am	happy	today	sad	feeling	sentiment
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{0}{4}$	$\frac{0}{4}$	1

$\frac{1}{4}$	$\frac{1}{4}$	$\frac{0}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{0}{4}$	-1
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{0}{4}$	$\frac{0}{4}$	$\frac{1}{4}$	1

The corresponding data frame having IDF values is :

i	am	happy	today	sad	feeling	sentiment
$\log(\frac{3}{3})$	$\log(\frac{3}{3})$	$\log(\frac{3}{2})$	$\log(\frac{3}{2})$	$\log(\frac{3}{1})$	$\log(\frac{3}{1})$	1
$\log(\frac{3}{3})$	$\log(\frac{3}{3})$	$\log(\frac{3}{2})$	$\log(\frac{3}{2})$	$\log(\frac{3}{1})$	$\log(\frac{3}{1})$	-1
$\log(\frac{3}{3})$	$\log(\frac{3}{3})$	$\log(\frac{3}{2})$	$\log(\frac{3}{2})$	$\log(\frac{3}{1})$	$\log(\frac{3}{1})$	1

The final dataframe with TF-IDF values is :

i	am	happy	today	sad	feeling	sentiment
0	0	0.044	0.044	0	0	1
0	0	0	0.044	0.119	0	-1
0	0	0.044	0	0	0.119	1

We can see that words like “happy”, “sad”, “feeling” have non-zero tf-idf value in the sentences in which they appear whereas “i”, “am” have zero tf-idf value. Intuitively also we can say that “sad” must be given more importance than “i” in determining the sentiment. TF-IDF gives more importance to words which occur frequently only in a few sentences and less importance to words which occur in a lot of sentences.

### Using algorithms for training and testing

Logistic Regression, Naive Bayes, Support Vector Machine (SVM), Random Forest, K-Nearest Neighbors were used for classifying independently. Ensemble learning techniques like XGBoost, Voting Classifier were also used.

Some of the algorithms used (for example Support Vector Machine) took more than a day for training and testing. Thus I used cuDF and cuML (since they can utilize GPU and thus can utilize google colab facility) in place of pandas and sklearn respectively for running large datasets with some of the algorithms.

Each algorithm was runned with different n-grams (i.e dataset with tweets having no tuples and different datasets having 2 to 4 tuples).

For every training and testing purpose, the dataset was divided into 80-20, respectively.

## Code and Output

Complete source code can be found on my github page ([link](#)).

Format of the confusion matrix printed in output is :

True Negatives	False Positives
False Negatives	True Positives

Negative refers to negative sentiment and Positive refers to positive sentiment.

Dataset was divided into 80-20 for training and testing every time.

Following are the screenshots of naive solution for different dataset sizes:

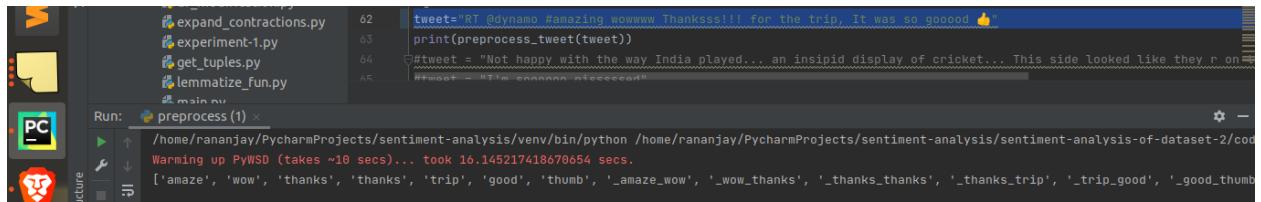
---

```
Confusion Matrix for LogisticRegression :  
[[18 34]  
 [ 6 41]]  
Accuracy :  
0.5959595959595959  
Confusion Matrix for SVM :  
[[25 27]  
 [ 9 38]]  
Accuracy :  
0.6363636363636364  
Total test data (around 20%) for both cases :  
99
```

---

```
Confusion Matrix for LogisticRegression :  
[[238 722]  
 [111 928]]  
Accuracy :  
0.5832916458229115  
Confusion Matrix for SVM :  
[[507 453]  
 [291 748]]  
Accuracy :  
0.6278139069534767  
Total test data (around 20%) for both cases :  
1999
```

Following is the screenshot of preprocessing of an example tweet :



A screenshot of the PyCharm IDE interface. On the left, there's a file tree with files like expand\_contractions.py, experiment-1.py, get\_tuples.py, lemmatize\_fun.py, and main.py. The main editor window shows Python code for preprocessing a tweet. The code includes importing libraries, defining a tweet variable, and running a preprocess function. A terminal window at the bottom shows the command run and its output.

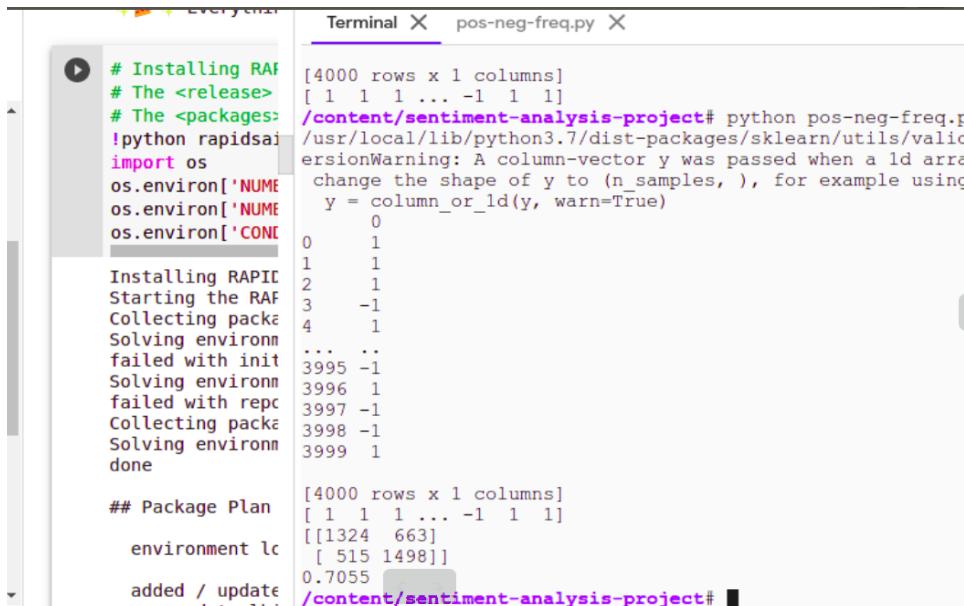
```
expand_contractions.py 62 tweet="RT @dynamo #amazing wowwww Thanksss!!! for the trip, It was so goodo 👍"
expand_contractions.py 63 print(preprocess_tweet(tweet))
expand_contractions.py 64 tweet = "Not happy with the way India played... an insipid display of cricket... This side looked like they r on
expand_contractions.py 65 ##wooo = "T@m sponono piecesend"
main.py
```

Run: preprocess(1) ×  
/home/rananjay/PycharmProjects/sentiment-analysis/venv/bin/python /home/rananjay/PycharmProjects/sentiment-analysis/sentiment-analysis-of-dataset-2/cod  
Warning: PyWSD (takes ~10 secs)... took 16.145217418678654 secs.  
['amaze', 'wow', 'thanks', 'thanks', 'trip', 'good', 'thumb', '\_amaze\_wow', '\_wow\_thanks', '\_thanks\_thanks', '\_thanks\_trip', '\_trip\_good', '\_good\_thumb']

Let's see the screenshots of all the experiments. Note that different screenshots have different train and test size (it can be seen from the sum of the values in the confusion matrix). This is because some models required more memory size and thus small dataset size was used for them. Evaluation of the output will be discussed in the next section.

Following are the screenshots of confusion matrix and accuracy obtained for dataset-1 and positive-negative frequency :

#### 1. positive-negative frequency, 1-gram, logistic regression.



A terminal window showing the execution of a Python script named pos-neg-freq.py. The script installs RAPID-SAI and runs a package plan. The output shows the execution of the script and its results, including a confusion matrix and accuracy.

```
# Installing RAF
# The <release>
# The <packages>
!python rapidssai
import os
os.environ['NUME
os.environ['NUME
os.environ['CONC

Installing RAPIC
Starting the RAF
Collecting packa
Solving environm
failed with init
Solving environm
failed with repc
Collecting packa
Solving environm
done

## Package Plan
environment lc
added / update
[4000 rows x 1 columns]
[ 1 1 1 ... -1 1 1]
[[1324 663]
 [ 515 1498]]
0.7055
```

#### 2. positive-negative frequency, 1-gram, SVM.

```

# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME'
os.environ['NUME'
os.environ['CONE'

Installing RAPIC
Starting the RAF
Collecting packa
Solving environm
failed with init

pos-neg-freq-svm.py
[content/sentiment-analysis-project]# python pos-neg-freq-svm.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:40: UserWarning: A column-vector y was passed when a 1d array was expected.
    "change the shape of y to (n_samples, ), for example using ravel()."
y = column_or_1d(y, warn=True)
[[1385  602]
 [ 628 1385]]
0.6925
[content/sentiment-analysis-project]#

```

### 3. positive-negative frequency, 1-gram, multinomial naive bayes.

```

# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME'
os.environ['NUME'
os.environ['CONE'

Installing RAPIC
Starting the RAF
Collecting packa
Solving environm
failed with init
Solving environm
failed with repc
Collecting packa
Solving environm
done
## Package Plan

[content/sentiment-analysis-project]# python pos-neg-freq-nb.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:40: UserWarning: A column-vector y was passed when a 1d array was expected.
    "change the shape of y to (n_samples, ), for example using ravel()."
y = column_or_1d(y, warn=True)
0
1
1
1
1
-1
1
...
3995 -1
3996 1
3997 -1
3998 -1
3999 1
[4000 rows x 1 columns]
[ 1  1  1 ... -1 -1  1]
[[1357  630]
 [ 561 1452]]
0.70225
[content/sentiment-analysis-project]#

```

### 4. positive-negative frequency, 1-gram, random forest.

```

# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME'
os.environ['NUME'
os.environ['CONE'

Installing RAPIC
Starting the RAF
Collecting packa
Solving environm
failed with init
Solving environm
failed with repc
Collecting packa
Solving environm
done
## Package Plan

environment lc

[content/sentiment-analysis-project]# python pos-neg-freq-rf.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:40: UserWarning: A column-vector y was passed when a 1d array was expected.
    "change the shape of y to (n_samples, ), for example using ravel()."
y = column_or_1d(y, warn=True)
3998 -1
3999 1
[content/sentiment-analysis-project]# python pos-neg-freq-rf.py:63: DataConversionWarning: A column-
a 1d array was expected. Please change the shape of y t
le using ravel().t-analysis-project#/content/sentim
model.fit(X_train,y_train)
0
1
1
1
1
-1
1
...
3995 -1
3996 1
3997 -1
3998 -1
3999 1
[4000 rows x 1 columns]
[ 1  1  1 ...  1  1 -1]
[content/sentiment-analysis-project]#
[ 689 1324]
0.6925
[content/sentiment-analysis-project]#

```

## 5. positive-negative frequency, 1-gram, KNN (K-Nearest Neighbors).

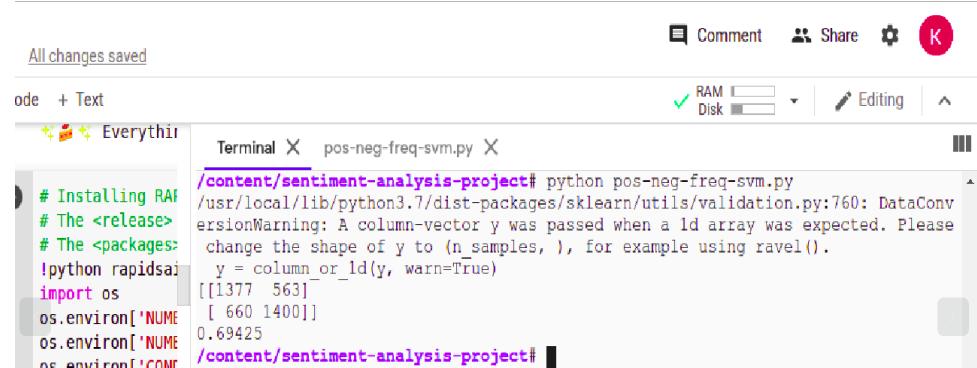
```
Terminal X pos-neg-freq-kn.py X
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python pos-neg-freq-kn.py
/usr/local/lib/python3.7/dist-packages/sklearn/neighbors/_classification.py:198: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y
(n_samples,), for example using ravel().
    return self._fit(X, y)
[[1449  538]
 [ 678 1335]]
0.696
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1#
```

## 6. positive-negative frequency, 2-grams, logistic regression.

```
# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUMEXPR_NUM_THREADS']=1
os.environ['NUMEXPR_MAX_THREADS']=1
os.environ['CONCURRENT_TASKS']=1

Installing RAPIDS
Starting the RAF
Collecting packages
Solving environment: failed with init
Solving environment: failed with repos
Collecting packages
Solving environment: done
[4000 rows x 1 columns]
[ 1  1  1 ...  1 -1  1]
## Package Plan
environment location: /content/sentiment-analysis-project
[ 543 1517]
0.70275
/content/sentiment-analysis-project#
```

## 7. positive-negative frequency, 2-grams, SVM.



All changes saved

ode + Text

RAM Disk

Comment Share

K

```
Terminal X pos-neg-freq-svm.py X
/content/sentiment-analysis-project# python pos-neg-freq-svm.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:760: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
    y = column_or_1d(y, warn=True)
[[1377  563]
 [ 660 1400]]
0.69425
/content/sentiment-analysis-project#
```

## 8. positive-negative frequency, 2-grams, multinomial naive bayes.

```

  [1]: ge the shape of y to (n samples, ), for example using ravel().
  [1]:     y = column_or_1d(y, warn=True)
  [1]: /content/sentiment-analysis-project# python pos-neg-freq-nb.py
  [1]: /usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:760: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n samples, ), for example using ravel().
  [1]:     y = column_or_1d(y, warn=True)
  [1]:         0
  [1]: 0    1
  [1]: 1    1
  [1]: 2    1
  [1]: 3    -1
  [1]: 4    1
  [1]: ...
  [1]: 3995  1
  [1]: 3996 -1
  [1]: 3997 -1
  [1]: 3998 -1
  [1]: 3999  1

  [4000 rows x 1 columns]
  [1] [ 1  1  1 ...  1 -1  1]
  [1] [[1326  614]
  [1] [ 585 1475]]
  [1] 0.70025

```

## 9. positive-negative frequency, 2-grams, random forest.

```

# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME
os.environ['NUME
Chromium Web Browser
os.environ['CONE
Installing RAPII
Starting the RAF
Collecting packa
Solving environm
failed with init
Solving environm
failed with repc
Collecting packa
Solving environm
done
## Package Plan
environment lc
pos-neg-freq-rf.py:63: DataConversionWarning
a 1d array was expected. Please change the s
le using ravel().
model.fit(X_train,y_train)
        0
0    1
1    1
2    1
3    -1
4    1
...
3995  1
3996 -1
3997 -1
3998 -1
3999  1

[4000 rows x 1 columns]
[1] [ 1  1  1 ...  1 -1  1]
[1] [[1428  512]
[1] [ 686 1374]]
0.7005
/content/sentiment-analysis-project#

```

## 10. positive-negative frequency, 2-grams, KNN.

```

All changes saved
ode + Text
Comment ✓ RAM Disk
* Everything Terminal × pos-neg-freq-kn.py ×
/content/sentiment-analysis-project# python pos-neg-freq-kn.py
pos-neg-freq-kn.py:63: DataConversionWarning: A column-vector y
a 1d array was expected. Please change the shape of y to (n_sam
ple using ravel().
model.fit(X_train,y_train)
[[1413  527]
 [ 696 1364]]
0.69425

```

## 11. positive-negative frequency, 3-grams, logistic regression.

```

# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME
os.environ['NUME
os.environ['CONC

Installing RAPII
Starting the RAF
Collecting packa
Solving environn
failed with init
Solving environn
failed with repc
Collecting packa
Solving environn
done
[4000 rows x 1 columns]
[ 1  1  1 ...  1 -1  1]
/content/sentiment-analysis-project# [ 543 1517]
0.70275
/content/sentiment-analysis-project# /content/sentim

```

## 12. positive-negative frequency, 3-grams, SVM.

```

# Installing RAF
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME
os.environ['NUME
os.environ['CONC

Installing RAPII
Starting the RAF
Collecting packa
Solving environn
failed with init
Solving environn
failed with repc
Collecting packa
Solving environn
done
[4000 rows x 1 columns]
[ 1  1  1 ...  1 -1  1]
/content/sentiment-analysis-project# python pos-neg-freq-svm.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:57: UserWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel(). This warning is displayed once per run.
y = column_or_1d(y, warn=True)
[[1377  563]
 [ 660 1400]]
0.69425
/content/sentiment-analysis-project#

```

## 13. positive-negative frequency, 3-grams, naive bayes.

```

# The <packages>
!python rapidsai
import os
os.environ['NUME
os.environ['NUME
os.environ['CONC

Installing RAPII
Starting the RAF
Collecting packa
Solving environn
failed with init
Solving environn
failed with repc
Collecting packa
Solving environn
done
[4000 rows x 1 columns]
[ 1  1  1 ...  1 -1  1]
[[1326  614]
 [ 585 1475]]
0.70025
added / updated

```

## 14. positive-negative frequency, 3-grams, random forest.

```

# Installing RAPID
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME']
os.environ['NUME']
os.environ['CONC']

Installing RAPID
Starting the RAF
Collecting packages
Solving environment
failed with init
Solving environment
failed with repc
Collecting packages
Solving environment
done
[[1428 512]
 [686 1374]]
## Package Plan
0.7005

```

15. positive-negative frequency, 3-grams, KNN.

```

# Installing RAPID
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUME']
os.environ['NUMF']

pos-neg-freq-kn.py:63: DataConversionWarning: A column-vector y
a 1d array was expected. Please change the shape of y to (n_samp
ple using ravel().
model.fit(X_train,y_train)
[[1413 527]
 [696 1364]]
0.69425

```

16. positive-negative frequency, 4-grams, logistic regression.

```

# The <packages>
!python rapidsai
import os
os.environ['NUME']
os.environ['NUME']
os.environ['CONC']

Installing RAPID
Starting the RAF
Collecting packages
Solving environment
failed with init
Solving environment
failed with repc
Collecting packages
Solving environment
done
[[1357 654]
 [473 1516]]
environment location
added / updated
0.71825

```

17. positive-negative frequency, 4-grams, SVM.

```
# Installing RAPID
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUMEXPR_NUM_THREADS'] = '1'
os.environ['NUMEXPR_MAX_THREADS'] = '1'
os.environ['CONCURRENT_FUTURE_THREADS'] = '1'

Installing RAPID
Starting the RAPID
Collecting packages
Solving environment: failed with initial attempt
Solving environment: failed with repodata
Collecting packages
```

```
Terminal × pos-neg-freq-svm.py × bag-of-words →
/content/sentiment-analysis-project# python pos-neg-freq-svm.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:65: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, 1) if you're passing column vectors. To avoid this warning, use ravel(). See https://scikit-learn.org/stable/user_guide/concepts.html#column-vector-warning for more information.
  y = column_or_1d(y, warn=True)
[[1434  577]
 [ 594 1395]]
0.70725
/content/sentiment-analysis-project#
```

#### 18. positive-negative frequency, 4-grams, naive bayes.

```
# Installing RAPID
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUMEXPR_NUM_THREADS'] = '1'
os.environ['NUMEXPR_MAX_THREADS'] = '1'
os.environ['CONCURRENT_FUTURE_THREADS'] = '1'

Installing RAPID
Starting the RAPID
Collecting packages
Solving environment: failed with initial attempt
Solving environment: failed with repodata
Collecting packages
Solving environment: done
## Package Plan
```

```
Terminal × pos-neg-freq.py ×
/content/sentiment-analysis-project# python pos-neg-freq.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:65: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, 1) if you're passing column vectors. To avoid this warning, use ravel(). See https://scikit-learn.org/stable/user_guide/concepts.html#column-vector-warning for more information.
  y = column_or_1d(y, warn=True)
      0
      1
      1
      1
      -1
      1
...
3995 -1
3996  1
3997 -1
3998  1
3999 -1
[4000 rows x 1 columns]
[ 1  1  1 ... -1  1 -1]
[[1396  615]
 [ 504 1485]]
0.72025
```

#### 19. positive-negative frequency, 4-grams, random forest.

```
# Installing RAPID
# The <release>
# The <packages>
!python rapidsai
import os
os.environ['NUMEXPR_NUM_THREADS'] = '1'
os.environ['NUMEXPR_MAX_THREADS'] = '1'
os.environ['CONCURRENT_FUTURE_THREADS'] = '1'

Installing RAPID
Starting the RAPID
Collecting packages
Solving environment: failed with initial attempt
Solving environment: failed with repodata
Collecting packages
Solving environment: done
## Package Plan
```

```
Terminal × pos-neg-freq.py ×
/content/sentiment-analysis-project# python pos-neg-freq.py
pos-neg-freq.py:65: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, 1) if you're passing column vectors. To avoid this warning, use ravel(). See https://scikit-learn.org/stable/user_guide/concepts.html#column-vector-warning for more information.
  model.fit(X_train,y_train)
      0
      1
      1
      1
      -1
      1
...
3995 -1
3996  1
3997 -1
3998  1
3999 -1
[4000 rows x 1 columns]
[ 1  1  1 ... -1  1 -1]
[[1520  491]
 [ 638 1351]]
0.71775
```

#### 20. positive-negative frequency, 4-grams, KNN.

```

COLLECTING PACKAGES
Solving environment
done
## Package Plan
environment located at /content/sentiment-analysis-project
added / updated
- cudatoolkit

```

```

/content/sentiment-analysis-project# python pos-neg-freq-kn.py
pos-neg-freq-kn.py:63: DataConversionWarning: A column-vector y was passed when
a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel(). -analysis-project#/content/sentim
    model.fit(X_train,y_train)
[[1518 493]
 [ 640 1349]]
0.71675
/content/sentiment-analysis-project#
[0] 0:python3*                                         "297a110bf8f4" 00:21 14-Nov-21

```

Following are the screenshots of confusion matrix and accuracy obtained for dataset-1 and bag of words :

### 1. bag of words, 1-gram, logistic regression.

```

# Installing RAPI
# The <release> command
# The <packages>
!python rapidsai-0.1.0-py3.7.egg
import os
os.environ['NUMBA_DISABLE_JIT'] = '1'
os.environ['NUMBA_NUM_THREADS'] = '1'
os.environ['CONDA_OVERRIDE_CUDA'] = '10.1'

Installing RAPIDS
Starting the RAPI
Collecting packages
Solving environment
failed with initia
Solving environment
failed with repod
Collecting packages
Solving environment
done

```

```

tcmalloc: large alloc 1216774144 bytes == 0x5639b9b980
1bb780354f 0x7f1bb7853b58 0x7f1bb7857b17 0x7f1bb78f620
d87240 0x5639b4dfb627 0x5639b4d88afa 0x5639b4dfad00 0;
da 0x5639b4df7737 0x5639b4df59ee 0x5639b4d88bda 0x5639
x5639b4df56f3 0x5639b4ebf4c2 0x5639b4ebf83d 0x5639b4e
9b4e96e0c 0x7f1bb938dbf7 0x5639b4e96cea
/usr/local/lib/python3.7/dist-packages/sklearn/linear_
convergenceWarning: lbgfs failed to converge (status=1)
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale
the features manually. Please also refer to the documentation for alternative
linear models
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)
[[1008 414]
 [ 357 1021]]
0.7246428571428571
/content/sentiment-analysis-project#

```

### 2. bag of words, 1-gram, SVM.

```

!apt update
import os
os._exit(0)

Updating your Colab environment
Hit:1 https://cloud.r-project.org
Ign:2 https://developer.downloadmirror.ubuntu.com
Hit:3 https://security.ubuntu.com
Ign:4 https://developer.downloadmirror.ubuntu.com
Hit:5 https://developer.downloadmirror.ubuntu.com
Hit:6 https://developer.downloadmirror.ubuntu.com
Hit:7 http://ppa.launchpad.net
Hit:8 http://archive.ubuntu.com

```

```

/tcmalloc: large alloc 1216774144 bytes == 0x557b14ca6000 @ 0x7f1
a001 0x7f154a15ae16 0x7f154a15aee6 0x7f154a19d585 0x7f154a1a97
54a233d65 0x557b108f8300 0x557b108f9ad8 0x557b109211d9 0x557b108
0x557b108f99c8 0x557b1092474a 0x557b10867af2 0x557b10896030 0x55
99c8 0x557b109211d9 0x557b10867af2 0x557b10896030 0x557b108f99c8
b1092474a 0x557b10867af2 0x557b10868d09 0x557b109438ab 0x557b109
0x557b109b1fd7 0x557b109b21ac 0x557b109b2709 0x557b109b285c 0x7f1
9bf7 0x557b1092f901
[[ 967 455]
 [ 309 1069]]
0.7271428571428571
/content/sentiment-analysis-project#

```

### 3. bag of words, 1-gram, naive bayes.

```

/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python bag-of-words-nb.py
/tcmalloc: large alloc 1216774144 bytes == 0x55f7df304000 @ 0x7fcfbefc4001 0x7fcfb24854f 0x7fcf
8 0x7fcfb29cb17 0x7fcfb33b203 0x55f7dbac4544 0x55f7dbac4240 0x55f7dbb38627 0x55f7dbac5afa 0x55
d00 0x55f7dbb329ee 0x55f7dbac5bda 0x55f7dbb34737 0x55f7dbb329ee 0x55f7dbac5bda 0x55f7dbb37d00 0x
329ee 0x55f7dbb326f3 0x55f7dbbf4c2 0x55f7dbbf4c2 0x55f7dbbf4c2 0x55f7dbbf4c2 0x55f7dbbf4c2 0x55f7dbbf4c2
d4acbf7 0x55f7dbbd3cea
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function
feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please
use get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
[[1065 357]
 [ 449 929]]
0.7121428571428572

```

### 4. bag of words, 1-gram, random forest.

I To view examples of instances, click the "Open Examples" button!

**OPEN EXAMPLES** **SEARCH STACK**

```
!ps -A
PID TTY      TIME
 1 ? 00:00:00
 8 ? 00:00:00
18 ? 00:00:00
47 ? 00:00:00
72 ? 00:00:00
73 ? 00:00:00
2371 ? 00:00:00
2391 ? 00:00:00
18352 ? 00:00:00
18353 pts/2 00:00:00
20266 ? 00:00:00
```

Terminal X cuda-cuml-test.py X bag-of-words-rf.py X

```
x5632cea659c8
/usr/local/lib/python3.7/site-packages/pandas/core/
WithCopyWarning:
A value is trying to be set on a copy of a slice fr
Try using .loc[row_indexer,col_indexer] = value ins
See the caveats in the documentation: https://panda
cs/stable/user_guide/indexing.html#returning-a-view
self[k1] = value[k2]
toward solving
/usr/local/lib/python3.7/site-packages/cuml/interna
67: UserWarning: To use pickling or GPU-based predi
g float32 data to fit the estimator
    ret_val = func(*args, **kwargs)
/usr/local/lib/python3.7/site-packages/cuml/interna
86: UserWarning: GPU based predict only accepts np.
l was trained on np.float64 data hence cannot use G
Defaulting to CPU-based Prediction.
To predict on float-64 data, set parameter predict_
    ret_val = func(*args, **kwargs)
<class 'cudf.core.series.Series'>
<class 'cudf.core.dataframe.DataFrame'>
we are here
[[ 921  501]
 [ 299 1079]]
0.7142857142857143
/content/sentiment-analysis-project#
```

## 5. bag of words, 1-gram, KNN.

```
!bash rapidsai-csp-utils/cc
import os
os._exit(00)

Updating your Colab environment
Hit:1 https://cloud.r-project.org
Ign:2 https://developer.download
Hit:3 http://security.ubuntu
Ign:4 https://developer.download
Hit:5 https://developer.download
Hit:6 https://developer.download
Hit:7 http://ppa.launchpad.
Hit:8 http://archive.ubuntu
Hit:11 http://archive.ubuntu
sdb / 0x0000499080901
/content/sentiment-analysis-project# python bag-of-words-kn.py
tcmalloc: large alloc 1216774144 bytes == 0x55d385178000 0 0x7f8e635b4a97 (0x63646d65 0x55d380b59300 0x55d380b5aad8 0x55d380b821d9 0x55d380af0x55d380b5a9c8 0x55d380b8574a 0x55d380ac8af2 0x55d380af7030 0x55d380b5a9c8 0x55d380b8574a 0x55d380ac8af2 0x55d380af7030 0x55d380b5a9c8 0x55d380c12fd7 0x55d380c131ac 0x55d380c13709 0x55d380c1385c 0x7f8e635b4a97 (0x842 580) [446 932]
0.6335714285714286
/content/sentiment-analysis-project# [0] 0:python+ "808a4101608d" 13:25 10-N
```

## 6. bag of words, 2-grams, logistic regression.

```
# Installing RAPIDS
# The <release> c
# The <packages>
!python rapidsai-
import os
os.environ['NUMBA_
os.environ['NUMBA_
os.environ['CONDA_
```

Installing RAPIDS  
Starting the RAPI  
Collecting package requirements...  
Solving environment: failed with initial attempt  
Solving environment: failed with repodata from https://rapidsai.github.io/rapidsai/conda/repodata.json  
Collecting package requirements...  
Solving environment: done  
  
## Package Plan ##  
environment location: /content/sentiment-analysis-project  
added / updated: 0.7510714285714286  
/content/sentiment-analysis-project#

## 7. bag of words, 2-grams, SVM.

```
psutil-5.8.0          low-level solving
psutil-5.8.0          <class 'cudf.core.series.Series'>
psutil-5.8.0          <class 'cudf.core.dataframe.DataFrame'>
we are here
[[ 898  409]
 [ 275 1018]]
0.7369230769230769
/content/sentiment-analysis-project#
```

## 8. bag of words, 2-grams, naive bayes.

```
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python bag-of-words-nb.py
tcmalloc: large alloc 5606662144 bytes == 0x56071e9d2000 @ 0x7fbf62884001 0x7fbf5f40854f 0x7f
8 0x7fbf5f45cb17 0x7fbf5f4fb203 0x560719edc544 0x560719edc240 0x560719f50627 0x560719eddafa 0x
d00 0x560719f4a9ee 0x560719eddafa 0x560719f4c737 0x560719f4a9ee 0x560719eddafa 0x560719f4fd00
4a9ee 0x560719f4a6f3 0x56071a0144c2 0x56071a01483d 0x56071a0146e6 0x560719fec163 0x560719febe0
166cbf7 0x560719febcea
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function
get_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. F
get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
tcmalloc: large alloc 44800 bytes == 0x56086ccc0000 @ 0x7fbf628821e7 0x7fbf5f40846e 0x7f
b 0x7fbf5f45935f 0x7fbf5f4fb103 0x560719edc544 0x560719edc240 0x560719f50627 0x560719f4a9ee 0x
bda 0x560719f4c737 0x560719f4a9ee 0x560719eddafa 0x560719f4c737 0x560719f4aced 0x560719eddafa
4b915 0x560719f4aced 0x560719eddafa 0x560719f4c737 0x560719f4a9ee 0x560719eddafa 0x560719f4c73
9f4a9ee 0x560719eddafa 0x560719f4c737 0x560719f4a9ee 0x560719eddafa 0x560719f4c737 0x560719f4a
719eddafa
[[844 268]
 [336 752]]
0.7254545454545455
```

## 9. bag of words, 2-grams, random forest.

```

import cudf
-----
ModuleNotFoundError:
<ipython-input-1>
ModuleNotFoundError:
-----
NOTE: If your
manually insta
To view exempl
"Open Examples
-----
OPEN EXAMPLES

```

## 10. bag of words, 2-grams, KNN.

```

Installing...
# you can now run the rest of th
import condacolab
condacolab.check()

Everything looks OK!

# Installing RAPIDS is now 'pyth
# The <release> options are 'sta
# The <packages> option are defa
!python rapidsai-csp-utils/colab
import os
os.environ['NUMBapro_NVVM'] = '/u
os.environ['NUMBapro_LIBDEVICE']
os.environ['CONDA_PREFIX'] = '/u

Installing RAPIDS Stable 21.10
Starting the RAPIDS install on C
Collecting package metadata (cur
Solving environment: ...working.
failed with initial frozen solve
Solving environment: ...working.
failed with repodata from curren
Collecting package metadata /ren

```

```

Terminal X bag-of-words-kn.py X
0x563532817d00 0x5635328129ee (
5328129ee 0x5635327a5bda 0x5635
26f3 0x5635328dc4c2 0x5635328dc
0x5635328b3e0c 0x7f60c6f1dbf7 (
tcmalloc: large alloc 525158809
0c81331e7 0x7f60c539346e 0x7f60
6103 0x5635327a4544 0x5635327a4
0x5635327a5bda 0x563532814737 (
532814737 0x563532812ced 0x5635
2ced 0x5635327a5bda 0x563532814
0x563532814737 0x5635328129ee (
5328129ee 0x5635327a5bda 0x5635
5bda
tcmalloc: large alloc 525158809
0c81331e7 0x7f60c539346e 0x7f60
3fe9 0x7f60c53e6d7d 0x7f60c53e1
0x5635327a4240 0x563532818627 (
088805cf 0x7f60888370e0 0x5635
4fe4 0x5635327a5afa 0x563532813
0x5635327a5afa 0x563532813915 (
5328dc4c2 0x5635328dc83d 0x5635
3e0c
/content/sentiment-analysis-pr
[333 835]
0.6295833333333334
/content/sentiment-analysis-pr

```

## 11. bag of words, 3-grams, logistic regression.

```

Solving envir
failed with r
Collecting pa
Solving envir
done
## Package Pl
environment
added / und

```

## 12. bag of words, 3-grams, SVM.

```

backports-1.0      | 4 KB  |         forward solving
backports-1.0      | 4 KB  | <class 'cudf.core.series.Series'>
<class 'cudf.core.dataframe.DataFrame'
we are here
jsonschema-4.2.1   | 118 K  | [ [595 322]
jsonschema-4.2.1   | 118 K  | [178 705]
jsonschema-4.2.1   | 118 K  | 0.7222222222222222
nonlocal-2.1 0.3.0 | 15.9  | /content/sentiment-analysis-project#

```

### 13. bag of words, 3-grams, naive bayes.

```

# Installing RA
# The <releas
# The <packag
!python rapid
import os
os.environ['N
os.environ['N
os.environ['C

Installing RA
Starting the
Collecting pa
Solving envir
failed with i
Solving envir
failed with r
Collecting pa
Solving envir
done

## Package Pl
environment
added / und
0.72

```

### 14. bag of words, 3-grams, random forest.

```

# Installing RA
# The <releas
# The <packag
!python rapid
import os
os.environ['N
os.environ['N
os.environ['C

Installing RA
Starting the
Collecting pa
Solving envir
failed with i
Solving envir
failed with r
Collecting pa
Solving envir
done

## Package Pl
environment
added / und
0.702

```

### 15. bag of words, 3-grams, KNN.

```

# Installing RAPIDS
# The <releases> options are
# The <packages> option are d
!python rapidsai-csp-utils/configure.py
import os
os.environ['NUMBapro_NVVM'] =
os.environ['NUMBapro_LIBDEVICE'] =
os.environ['CONDA_PREFIX'] = 

Installing RAPIDS Stable 21.1
Starting the RAPIDS install
Collecting package metadata (Solving environment: ...working)
failed with initial frozen so
Solving environment: ...working
failed with repodata from cur
Collecting package metadata (Solving environment: ...working)
done

## Package Plan ##

environment location: /usr/

```

/content/sentiment-analysis-project# [437 446]

0.5972222222222222

/content/sentiment-analysis-project#

## 16. bag of words, 4-grams, logistic regression.

```

# Installing RAPIDS is now 'p'
# The <release> options are
# The <packages> option are d
!python rapidsai-csp-utils/configure.py
import os
os.environ['NUMBapro_NVVM'] =
os.environ['NUMBapro_LIBDEVICE'] =
os.environ['CONDA_PREFIX'] = 

Installing RAPIDS Stable 21.1
Starting the RAPIDS install
Collecting package metadata (Solving environment: ...working)
failed with initial frozen so
Solving environment: ...working
failed with repodata from cur
Collecting package metadata (Solving environment: ...working)
done

## Package Plan ##

environment location: /usr/

```

## 17. bag of words, 4-grams, SVM.

The screenshot shows a Jupyter Notebook interface with two code cells and a terminal output.

```
# Installing RAPIDS is
# The <release> option:
# The <packages> option:
!python rapidsai-csp-u
import os
os.environ['NUMBapro_N'
os.environ['NUMBapro_L'
os.environ['CONDA_PREF

protobuf-3.16.0
protobuf-3.16.0
protobuf-3.16.0
google-resumable-med
google-resumable-med
llvmlite-0.36.0
llvmlite-0.36.0
llvmlite-0.36.0
```

Terminal output:

```
c8 0x55cff06764ac
0x55cff05bdaf2 0x55cff05ec030 0
0x55cff064f9c8 0x55cff06764ac
8 0x55cff06771d9 0x55cff05bdaf2
d9 0x55cff05bdaf2
0x55cff05ec030 0x55cff064f9c8 0
/usr/local/lib/python3.7/site-p
ithCopyWarning:
A value is trying to be set on
Try using .loc[row_indexer,col_
protobuf-3.16.0 | See the caveats in the documen
protobuf-3.16.0 | s/stable/user_guide/indexing.ht
protobuf-3.16.0 | self[k1] = value[k2]
google-resumable-med | toward solving
google-resumable-med | <class 'cudf.core.series.Series'
we are here | <class 'cudf.core.dataframe.Dat
llvmlite-0.36.0 | we are here
llvmlite-0.36.0 | [[523 293]
llvmlite-0.36.0 | [153 631]]
llvmlite-0.36.0 | 0.72125
/ccontent/sentiment-analysis-p
```

## 18. bag of words, 4-grams, naive bayes.

The screenshot shows a Jupyter Notebook interface with a code cell and a terminal output.

```
-
```

```
ModuleNotFoundError
last)
<ipython-input-1-e13365c50bc4> in
----> 1 import cudf

ModuleNotFoundError: No module nam
-
```

**NOTE: If your import is failing due to a missing dependency, please manually install dependencies using pip or conda.**

```
To view examples of installing some packages, click the "Open Examples" button below.
```

OPEN EXAMPLES SEARCH STACK OVERFLOW

Terminal output:

```
3d1f1d79 0x7f803d1f4e4c 0x7f80
1be3d 0x7f803d31d516 0x558d917c3
db 0x558d9183d0b2 0x558d917c3
0x558d917c3fe4 0x558d917c19ee
58d91753437 0x558d91753240 0x
91754bda 0x558d917c3737 0x558
c6d00 0x558d917c19ee
tcmalloc: large alloc 3251347
f803fe851e7 0x7f803d0e546e 0x
3d1c83a9 0x7f803d1caab5 0x558
c19ee 0x558d91754bda 0x558d917
da 0x558d917c2c0d 0x558d91754
0x558d917c16f3 0x558d9188b4c2
58d91863163 0x558d91862e0c 0x
tcmalloc: large alloc 3251347
f803fe851e7 0x7f803d0e546e 0x
3d1f1d79 0x7f803d1f4e4c 0x7f80
1be3d 0x7f803d31d516 0x558d917
db 0x558d9183d0b2 0x558d917c3
0x558d917c3fe4 0x558d917c19ee
58d91753437 0x558d91753240 0x
91754bda 0x558d917c3737 0x558
c6d00 0x558d917c19ee
[[325 491]
 [174 610]]
0.584375
/ccontent/sentiment-analysis-p
```

## 19. bag of words, 4-grams, random forest.

```

ModuleNotFoundError
last)
<ipython-input-1-e13365c50bc4> in <module>
----> 1 import cudf

ModuleNotFoundError: No module named 'cu

NOTE: If your import is failing due to a
manually install dependencies using eith

To view examples of installing some comm
"Open Examples" button below.

OPEN EXAMPLES SEARCH STACK OVERFLOW

```

x5585fb0159ee 0x5585fafafa8bda  
0x5585fafafa8bda 0x5585fb01773  
da  
tcmalloc: large alloc 1612070  
0x7f33a33661e7 0x7f33a05c64  
35f 0x7f33a06b9103 0x5585fafafa  
1b627 0x5585fb0159ee 0x5585fa  
b0159ee 0x5585fafafa8bda 0x5585  
5fafafa8bda 0x5585fb016915 0x55  
85fb017737 0x5585fb0159ee 0x  
x5585fb0159ee 0x5585fafafa8bda  
0x5585fafafa8bda 0x5585fb01773  
da  
tcmalloc: large alloc 3224141  
0x7f33a33661e7 0x7f33a05c64  
d97 0x7f33a0616fe9 0x7f33a061  
bb761 0x5585fafafa7544 0x5585fa  
b0159ee 0x5585fafafa8bda 0x5585  
5fafafa9271 0x5585fafafa775f 0x55  
f33a0618a93 0x7f33a06190bc 0x  
x7f33a06bb761 0x5585fafafa7544  
0x5585fb0159ee 0x5585fafafa8b  
ee  
[[383 433]  
[ 70 714]]  
0.685625

## 20. bag of words, 4-grams, KNN.

```

cols = y.columns
y[cols] = y[cols].apply(pd.to_numeric, e

x = cudf.DataFrame.from_pandas(x).astype
y = cudf.DataFrame.from_pandas(y).astype

# x=x.astype('float64')
# y=y.astype('float64')
# print(x.head(2))

print("toward solving")
solve(x, y)

import cudf

ModuleNotFoundError
last)
<ipython-input-1-e13365c50bc4> in <module>
----> 1 import cudf

ModuleNotFoundError: No module named 'cu

```

11737 0x55db1ac0f9ee 0x55db1a1  
1ac0f9ee 0x55db1aba2bda 0x55d  
5db1aba2bda  
tcmalloc: large alloc 1612070@  
0x7fbf926251e7 0x7fbf8f8854  
d635f 0x7fbf8f978103 0x55db1a1  
1ac15627 0x55db1ac0f9ee 0x55d  
5db1ac0f9ee 0x55db1aba2bda 0x  
0x55db1aba2bda 0x55db1ac10915  
da 0x55db1ac11737 0x55db1ac0f9  
11737 0x55db1ac0f9ee 0x55db1a1  
1ac0f9ee 0x55db1aba2bda 0x55d  
5db1aba2bda  
tcmalloc: large alloc 6448283@  
0x7fbf926251e7 0x7fbf8f8854  
d5d97 0x7fbf8f8d5fe9 0x7fbf8f8  
8f97a761 0x55db1aba1544 0x55d  
5db1ac0f9ee 0x55db1aae1e2b 0x  
0x55db1aae1e2b 0x55db1aba3698  
fa 0x55db1ac10c0d 0x55db1aba2  
a2afa 0x55db1ac10915 0x55db1a  
1acd94c2 0x55db1acd983d 0x55d  
5db1acb0e0c  
[[624 192]  
[445 339]]  
0.601875

Following are the screenshots of confusion matrix and accuracy obtained for dataset-1 and TF-IDF :

## 1. TF-IDF, 1-gram, logistic regression.

```

import os
os.environ['N
os.environ['N
os.environ['C

Installing RA
Starting the
Collecting pa
Solving envir
failed with i
Solving envir

/content/sentiment-analysis-project# python tf-idf-lr.py
tcmalloc: large alloc 1216774144 bytes == 0x561089314000 @
09de16 0x7fcf7809dee6 0x7fcf780e0585 0x7fcf780e4a97 0x7fcf
x561084f77ad8 0x561084f9fld9 0x561084f13e94 0x561084f779c8
e5af2 0x561084f14030 0x561084f779c8 0x561084f9fld9 0x56108
561084f779c8 0x561084fa274a 0x561084ee5af2 0x561084ee6d09
5f53 0x56108502ffd7 0x5610850301ac 0x561085030709 0x56108
61084fad901
[[1012 410]
 [ 331 1047]]
0.7353571428571428

```

## 2. TF-IDF, 1-gram, SVM.

[ To view examples of instances ]

"Open Examples" button | -----

-

OPEN EXAMPLES SEARCH STACK

! ps -A

PID	TTY	TII
1 ?	00:00:	
8 ?	00:00:	
18 ?	00:00:	
47 ?	00:00:	
72 ?	00:00:	
73 ?	00:00:	
2371 ?	00:00:	
2391 ?	00:00:	
18352 ?	00:00:	
18353 pts/2	00:00:	
20266 ?	00:00:	

Terminal X cuda-cuml-test.py :

See the caveats in the documentation: https://cuml.apache.org/stable/user\_guide/indexing.html#returning-views

```
self._setitem_single_block(indexer, value, tcmalloc: large alloc 121677: 7 0x7f66736dfd0d 0x7f66736df0b85 0x563631216300 0x5636312179c8 0x563631216300 0x5636312179c8 0x5636312179c8 /usr/local/lib/python3.7/site-packages/cuml/test/test_cuml.py:105: WithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value or .at[row_indexer,col_indexer] = value. See the caveats in the documentation: https://cuml.apache.org/stable/user_guide/indexing.html#returning-views self[k1] = value[k2]
```

toward solving <class 'cudf.core.series.Series'>

<class 'cudf.core.dataframe.DataFrame'>

we are here

```
[[ 996 4261]
 [ 326 1052]]
0.7314285714285714
```

## 3. TF-IDF, 1-gram, random forest.

Terminal X tf-idf-rf.py :

/content/sentiment-anal

```
/usr/local/lib/python3.7/site-packages/pandas/core/generic.py:103: UserWarning: on get_feature_names is deprecated in 1.2. Please use get_feature_names instead.
warnings.warn(msg, category=UserWarning)
[[471 435]
 [123 771]]
0.69
/content/sentiment-anal
```

## 4. TF-IDF, 1-gram, KNN.

Q F D

..

sa\_file\_processed.csv

sa\_file\_processed1.csv

sa\_file\_processed2.csv

sa\_file\_processed\_three\_tupl...

sentiment2.csv

spell-correction.py

spell\_correct.py

test.py

test1.py

tf-idf-cuda.py

tf-idf-gnb.py

tf-idf-kn.py

tf-idf-lr.py

tf-idf-nb.py

tf-idf-rf.py

tf-idf-svm.py

condacolab\_install.log

Disk 113.00 GB available

[ To view examples of instances ]

"Open Examples" button | -----

-

OPEN EXAMPLES SEARCH STACK

! ps -A

PID	TTY	TII
1 ?	00:00:	
8 ?	00:00:	
18 ?	00:00:	
47 ?	00:00:	
72 ?	00:00:	
73 ?	00:00:	
2371 ?	00:00:	
2391 ?	00:00:	
18352 ?	00:00:	
18353 pts/2	00:00:	
20266 ?	00:00:	

Terminal X cuda-cuml-test.py X tf-idf-kn.py X

See the caveats in the documentation: https://cuml.apache.org/stable/user\_guide/indexing.html#returning-views

```
self._setitem_single_block(indexer, value, tcmalloc: large alloc 1216774144 bytes == 0x7f66736fd0b85 0x7f66736fd0b85 0x56207c027300 0x56207c028ad8 0x56207c05b85 0x56207c04f4ac 0x56207bf96af2 0x56207bfc5030 0x56207bf96af2 0x56207bfc5030 0x56207c0289c8 0x56207c0289c8 /usr/local/lib/python3.7/site-packages/pandas/core/generic.py:103: WithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value or .at[row_indexer,col_indexer] = value. See the caveats in the documentation: https://cuml.apache.org/stable/user_guide/indexing.html#returning-views self[k1] = value[k2]
```

toward solving <class 'cudf.core.series.Series'>

<class 'cudf.core.dataframe.DataFrame'>

we are here

```
[[ 282 1140]
 [ 107 1271]]
0.5546428571428571
/content/sentiment-analysis-project#
```

## 5. TF-IDF, 2-grams, logistic regression.

```

----- /content/sentiment-analysis-project# python tf-idf-lr.py
tcmalloc: large alloc 8728051712 bytes == 0x55da88078000
2d6e16 0x7ff04c2d6ee6 0x7ff04c319585 0x7ff04c31da97 0x7ff
x55da84492ad8 0x55da844ba1d9 0x55da8442ee94 0x55da844929c
00af2 0x55da8442f030 0x55da844929c8 0x55da844ba1d9 0x55da
55da844929c8 0x55da844bd74a 0x55da84400af2 0x55da84401d09
0f53 0x55da8454afd7 0x55da8454blac 0x55da8454bd09 0x55da8
5da844c8901
tcmalloc: large alloc 6982443008 bytes == 0x55dc90c4a000
2d6d0d 0x7ff04c2d6d87 0x7ff04c319574 0x7ff04c31a0ff 0x7ff
x55da84492ad8 0x55da844ba1d9 0x55da8442ee94 0x55da844929c
00af2 0x55da8442f030 0x55da844929c8 0x55da844ba1d9 0x55da
55da844929c8 0x55da844ba1d9 0x55da84400af2 0x55da8442f030
94ac 0x55da84400af2 0x55da8442f030 0x55da844929c8 0x55da8
5da8442f030 0x55da844929c8
tcmalloc: large alloc 1745616896 bytes == 0x55de30f44000
2d6d0d 0x7ff04c2d6d87 0x7ff04c319574 0x7ff04c31a0ff 0x7ff
x55da84492ad8 0x55da844ba1d9 0x55da8442ee94 0x55da844929c
00af2 0x55da8442f030 0x55da844929c8 0x55da844ba1d9 0x55da
55da844929c8 0x55da844ba1d9 0x55da84400af2 0x55da8442f030
94ac 0x55da84400af2 0x55da8442f030 0x55da844929c8 0x55da8
5da8442f030 0x55da844929c8
-----
```

## 6. TF-IDF, 2-grams, SVM.

```

----- [1] Done in 0.0s.
[2] Restarting
----- [2] # you can now
import condacolab
condacolab.ch
+ Everyt

# Installing RA
Starting the
Collecting pa
Solving envir
failed with i
Solving envir 0.7396153846153846
-----
```

Terminal X cuda-cuml-test.py X tf-idf-svm.py X

```

aaaf2 0x55c6970d9030 0x55c69713c9c8 0x55c69716
55c69713c9c8 0x55c6971641d9 0x55c6970aaaaf2 0x5
41d9 0x55c6970aaaaf2 0x55c6970d9030 0x55c69713c
5c6970d9030 0x55c69713c9c8
tcmalloc: large alloc 7610204160 bytes == 0x55
b73d0d 0x7f388db73d87 0x7f388dbb6574 0x7f388db
x55c69713cad8 0x55c6971641d9 0x55c6970d8e94 0x
aaaaf2 0x55c6970d9030 0x55c69713c9c8 0x55c69716
55c6971634ac 0x55c6970aaaaf2 0x55c6970d9030 0x5
aaaf2 0x55c6970d9030 0x55c69713c9c8 0x55c697164
5c69713c9c8 0x55c6971641d9
/usr/local/lib/python3.7/site-packages/pandas/rnning:
A value is trying to be set on a copy of a sli
Try using .loc[row_indexer,col_indexer] = valu
See the caveats in the documentation: https://u
user_guide/indexing.html#returning-a-view-vers
self[k1] = value[k2]
toward solving
<class 'cudf.core.series.Series'>
<class 'cudf.core.dataframe.DataFrame'>
we are here
[[ 917 390]
 [ 287 1006]]
0.7396153846153846
-----
```

## 7. TF-IDF, 2-grams, random forest.

```

----- /content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python tf-idf-rf.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Func
et feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed
. Please use get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
[[476 430]
 [135 759]]
0.6861111111111111
-----
```

## 8. TF-IDF, 2-grams, KNN.

```

----- /content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python tf-idf-kn.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Fu
ture_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in
use get_feature_names_out instead./sentiment-analysis-of-dataset-1# python tf-idf-kn.py
warnings.warn(msg, category=FutureWarning)rn/utils/deprecation.py:87: FutureWarning: Fu
[[160 746]]entiment-analysis-project/sentiment-analysis-of-dataset-1# d will be removed in
[ 21 873]
0.5738888888888889
-----
```

## 9. TF-IDF, 3-grams, logistic regression.

```
# Installing RAPIDS is now 'p
# The <release> options are '
# The <packages> option are c
!python rapidsai-csp-utils/cc
import os
os.environ['NUMBAPRO_NVVM'] =
os.environ['NUMBAPRO_LIBDEVIC
os.environ['CONDA_PREFIX'] =
Installing RAPIDS Stable 21.1
Starting the RAPIDS install o
Collecting package metadata (
Solving environment: ...worki
failed with initial frozen so
Solving environment: ...worki
failed with repodata from cur
Collecting package metadata (
Solving environment: ...worki
done
## Package Plan ##
environment location: /usr/
added / updated speci.
0.729
```

```
Terminal X tf-idf-lr.py
185001 0x7f52b03c5e16 0
x7f52b049ed65 0x5584187
e8e94 0x55841874c9c8 0x
55841874c9c8 0x55841877
c9c8 0x55841877774a 0x5
584187faf53 0x558418804
85c 0x7f52b11b4bf7 0x55
tcmalloc: large alloc 7
1831e7 0x7f52b03c5d0d 0
x7f52b049eb85 0x5584187
e8e94 0x55841874c9c8 0x
55841874c9c8 0x55841877
c9c8 0x5584187741d9 0x5
584187734ac 0x5584186ba
1d9 0x5584186baaf2 0x55
tcmalloc: large alloc 1
1831e7 0x7f52b03c5d0d 0
x7f52b049eb85 0x5584187
e8e94 0x55841874c9c8 0x
55841874c9c8 0x55841877
c9c8 0x5584187741d9 0x5
584187734ac 0x5584186ba
1d9 0x5584186baaf2 0x55
[[690 323]
[219 768]]
0.729
```

## 10. TF-IDF, 3-grams, SVM.

```
# Installing RAPIDS is now 'p
# The <release> options are '
# The <packages> option are c
!python rapidsai-csp-utils/cc
import os
os.environ['NUMBAPRO_NVVM'] =
os.environ['NUMBAPRO_LIBDEVIC
os.environ['CONDA_PREFIX'] =
Installing RAPIDS Stable 21.1
Starting the RAPIDS install o
Collecting package metadata (
Solving environment: ...worki
failed with initial frozen so
Solving environment: ...worki
failed with repodata from cur
Collecting package metadata (
Solving environment: ...worki
done
## Package Plan ##
environment location: /usr/
added / updated speci.
```

```
Terminal X cuda-cuml-test.py X tf-idf-svm.py
573fd0fc1d9 0x5573fd042af2 0x5573fd071030 0
1d9 0x5573fd042af2 0x5573fd071030 0x5573fd0
tcmalloc: large alloc 7232405504 bytes == 0
51ble7 0x7fee9271dd0d 0x7fee9271dd87 0x7fee
x7fee92756b85 0x5573fd0d3300 0x5573fd0d4ad8
70e94 0x5573fd0d49c8 0x5573fd0fb4ac 0x5573f
5573fd0d49c8 0x5573fd0fc1d9 0x5573fd070e94
b4ac 0x5573fd042af2 0x5573fd071030 0x5573fd
573fd042af2 0x5573fd071030 0x5573fd0d49c8 0
af2 0x5573fd071030 0x5573fd0d49c8 0x5573fd0
/usr/local/lib/python3.7/site-packages/pand
ttingWithCopyWarning:
A value is trying to be set on a copy of a
Try using .loc[row_indexer,col_indexer] = v
See the caveats in the documentation: https
as-docs/stable/user_guide/indexing.html#ret
py
    self[k1] = value[k2]
toward solving
<class 'cudf.core.series.Series'>
<class 'cudf.core.dataframe.DataFrame'>
we are here
[[599 318]
[164 719]]
0.7322222222222222
```

## 11. TF-IDF, 3-grams, random forest.

manually install dependencies using [pip](#)

To view examples of installing some common dependencies, click the "Open Examples" button below.

```
!ps -A
```

PID	TTY	TIME	CMD
1	?	00:00:00	docker-in
8	?	00:00:08	node
18	?	00:00:00	tail
47	?	00:00:03	colab-finder
72	?	00:00:08	jupyter-lab
73	?	00:00:03	dap_middleware
2371	?	00:00:22	python3
2391	?	00:00:02	python3
18352	?	00:00:00	tmux: session
18353	pts/2	00:00:00	bash
20266	?	00:00:00	ps

```
Terminal X tf-idf-rf.py : da 0x55b73fb70737 0x55b73fb6ec0 0x55b73fb6ec0 0x55b73fb70737 0x55b73fb6e9ee 0x55b73fb01bda 0x55b73fb70737 0x55b73fb6e9ee 0x55b73fb01bda tcmalloc: large alloc 1 f4a215f61e7 0x7f4ale856 1e949103 0x55b73fb00544 0x55b73fb6ec0 0x55b73fb01bda da 0x55b73fb70737 0x55b73fb6ec0 0x55b73fb6ec0 0x55b73fb70737 0x55b73fb6e9ee 0x55b73fb01bda tcmalloc: large alloc 1 f4a215f61e7 0x7f4ale856 1e949103 0x55b73fb00544 0x55b73fb6ec0 0x55b73fb01bda tcmalloc: large alloc 3 f4a215f61e7 0x7f4ale856 1e8a6fe9 0x7f4ale8a9d7d 00544 0x55b73fb00240 0x55b73fb70737 0x55b73fb6ec0 0x55b73fb03305 0x7f4ale8a9cbb 3fb00240 0x55b73fb74627 70737 0x55b73fb6e9ee [[481 532] [102 885]] 0.683 /content/sentiment-analysis
```

## 12. TF-IDF, 3-grams, KNN.

```
Terminal X tf-idf-knn.py : 8b1bed00 0x56508b1b99ee 3 0x56508b25ae0c 0x7f5f5f /usr/local/lib/python3.7/site-packages/sklearn/feature_selection/_chfscore.py:140: DeprecationWarning: Feature selection module is deprecated, Please use get_feature_names instead. warnings.warn(msg, category=DeprecationWarning) tcmalloc: large alloc 5 /content/sentiment-analysis 0x561f4e3cd2d00 0x561f4e3c2d00 0x561f4e3bd9ee 3 0x561f4e45ee0c 0x7fa5f /usr/local/lib/python3.7/site-packages/sklearn/feature_selection/_chfscore.py:140: DeprecationWarning: Feature selection module is deprecated, Please use get_feature_names instead. warnings.warn(msg, category=DeprecationWarning) tcmalloc: large alloc 4 fa80c7b 0x7fa4ffa8135f 0x561f4e350bda 0x561f4e350bda 0x561f4e3be915 4e350bda 0x561f4e3be915 0x561f4e3bf737 0x561f4e3bf737 0x561f4e3bdce 0x561f4e3bdce [[506 306] [163 625]] 0.706875
```

## 13. TF-IDF, 4-grams, logistic regression.

```

tcmalloc: large alloc 8446203040 bytes == 0x55511f1dec000 @ 0x7f04c3740001 0x7f04c02c454f
[content/sentiment-analysis-project/sentiment-analysis-of-dataset-1]# python tf-idf-lr.py
tcmalloc: large alloc 6217400320 bytes == 0x5557974e0000 @ 0x7f04c3740001 0x7f04c02c454f
0314b58 0x7f04c0318b17 0x7f04c03b7203 0x5557938cd544 0x5557938cd240 0x555793941627 0x5557
0x555793940d00 0x55579393b9ee 0x5557938cebd 0x55579393d737 0x55579393b9ee 0x5557938cebd
93940d00 0x55579393b9ee 0x55579393b6f3 0x55579393a054c2 0x555793a0583d 0x555793a056e6 0x5557
0x55579393dce0c 0x7f04c2528bf7 0x55579393dcea
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Fu
et_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be remove
. Please use get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
tcmalloc: large alloc 4973920256 bytes == 0x55590a040000 @ 0x7f04c373e1e7 0x7f04c02c446e
0314c7b 0x7f04c031535f 0x7f04c03b7103 0x5557938cd544 0x5557938cd240 0x555793941627 0x5557
0x5557938cebd 0x55579393d737 0x55579393b9ee 0x5557938cebd 0x55579393d737 0x55579393bce
938cebd 0x55579393c915 0x55579393bcfd 0x5557938cebd 0x55579393d737 0x55579393b9ee 0x5557
a 0x55579393d737 0x55579393b9ee 0x5557938cebd 0x55579393d737 0x55579393b9ee 0x5557938ce
79393d737 0x55579393bcfd 0x5557938cebd
[[479 219]
 [161 541]]
0.7285714285714285

```

#### 14. TF-IDF, 4-grams, SVM.

```

] ⏺ Done in 0:00:23
🔁 Restarting kernel...
[1] # you can now run the rest of
    import condacolab
    condacolab.check()

    ✨ Everything looks OK!

    # Installing RAPIDS is now 'p
    # The <release> options are '
    # The <packages> option are d
!python rapidsai-csp-utils/co
import os
os.environ['NUMBapro_NVVM'] =
os.environ['NUMBapro_LIBDEVIC
os.environ['CONDA_PREFIX'] =
Installing RAPIDS Stable 21.1
Starting the RAPIDS install o
Collecting package metadata (
Solving environment: ...worki
failed with initial frozen so
Solving environment: ...worki

```

```

Terminal ✘ cuda-cuml-test.py
65296c761d9 0x565296bbcaf2 0i
1d9 0x565296bbcaf2 0x565296be
tcmalloc: large alloc 8060354
ddf1e7 0x7f3d3f021d0d 0x7f3d3
x7f3d3f0fab85 0x565296c4d300
eae94 0x565296c4e9c8 0x565296
565296c4e9c8 0x565296c761d9 (
54ac 0x565296bbcaf2 0x5652961
65296bbcaf2 0x565296beb030 0i
af2 0x565296beb030 0x565296c4
/usr/local/lib/python3.7/site-
ttingWithCopyWarning:
A value is trying to be set c
Try using .loc[row_indexer,co
See the caveats in the docume
as-docs/stable/user_guide/inst
py
    self[k1] = value[k2]
toward solving
<class 'cudf.core.series.Series'
<class 'cudf.core.dataframe.DataFrame'
we are here
[[544 272]
 [169 615]]
0.724375

```

#### 15. TF-IDF, 4-grams, random forest.

```

] ⏺ Done in 0:
🔁 Restarting kernel...
[2] # you can now
    import condacolab
    condacolab.ch

    ✨ Everything looks OK!

    # Installing RA
    # The <releas
    # The <packag
!python rapid
import os
os.environ['N
os.environ['N
os.environ['C
Installing RA
Starting the
Collecting pa
Solving envir

```

```

Terminal ✘ cuda-cuml-test.py
/usr/local/lib/python3.7/site-p
rning: To use pickling or GPU-ba
t the estimator
    ret_val = func(*args, **kwargs
/usr/local/lib/python3.7/site-p
rning: GPU based predict only ac
float64 data hence cannot use GF
Defaulting to CPU-based Predicti
To predict on float-64 data, set
    ret_val = func(*args, **kwargs
tcmalloc: large alloc 1612070912
62f136 0x7f2d23610bea 0x7f2d6ab9
x55ddcd8b63a8 0x55ddcd91bf93 0x7
6ea17 0x55ddcd8b4af2 0x55ddcd8e3
55ddcd9469c8 0x55ddcd96d4ac 0x55
4f53 0x55ddcd9fefd7 0x55ddcd9ff1
5ddcd97c901
<class 'cudf.core.series.Series'
<class 'cudf.core.dataframe.DataFrame'
we are here
[content/sentiment-analysis-proj
[ 73 711]
0.6775
[content/sentiment-analysis-proj

```

## 16. TF-IDF, 4-grams, KNN.

```
Terminal X tf-idf-kn.py X
f98abd00 0x563ff98a69ee 0x
3 0x563ff9947e0c 0x7ff9775
/usr/local/lib/python3.7/d
et_feature_names is deprec
. Please use get_feature_n
warnings.warn(msg, catég
tcmalloc: large alloc 6448
/content/sentiment-analyti
tcmalloc: large alloc 6217
5daeb58 0x7f5ed5db2b17 0x7
0x559ecf185d00 0x559ecf18
cf185d00 0x559ecf1809ee 0x
3 0x559ecf221e0c 0x7f5ed7f
/usr/local/lib/python3.7/d
et_feature_names is deprec
. Please use get_feature_n
warnings.warn(msg, catég
tcmalloc: large alloc 4973
5daec7b 0x7f5ed5daf35f 0x7
0x559ecf113bda 0x559ecf18
cf113bda 0x559ecf181915 0x
a 0x559ecf182737 0x559ecf1
ecf182737 0x559ecf180ced 0
[[438 260]
[143 559]]
0.7121428571428572
-----
```

Following are the screenshots of confusion matrix and accuracy obtained for dataset-1 and XGBoost :

### 1. bag of words, 1-gram.

```
Building dependency 1
Reading state information...
Selected version '11.3.0-17ubuntu1'
The following additional packages will be installed:
  gcc-11-base libgcc-11-dev libstdc++-11-dev
The following NEW packages will be installed:
  gcc-11-base libgcc-11-dev libstdc++-11-dev
The following packages will be upgraded:
  libstdc++6
1 upgraded, 3 newly installed, 0 to remove and 0 not upgraded.
# This will install (3) new packages.
import condacolab
/stable/user_guide/index.html
[ 0 :bash*
```

### 2. bag of words, 2-grams.

```
gcc-11-base libgcc-11-dev libstdc++-11-dev
The following packages will be installed:
libstdc++6
1 upgraded, 3 newly installed, 0 to remove and 0 not upgraded.
# This will install (3) new packages.
import condacolab
/tcmalloc: large alloc 376
x7fb3e87413b5 0x7fb3e8741
7fb3e8602c95 0x7fb3e8615f
fb411cf65bc 0x7fb411cf59e
60d9975bda 0x5560d99e3c0c
0d99e29ee 0x5560d9975bda
d99e3c0d 0x5560d9975afa ([560 494]
[135 811])
0.6855
-----
```

### 3. bag of words, 4-grams.

```
the following packages will be upgraded:
  libstdc++6 7fe93e654c95 0x7f
  1 upgraded, 0 newly installed, 0 to remove and 0 not upgraded.
# This will install 0.665
```

Hard voting was used for ensemble learning. I chose algorithms giving the highest accuracy and then final sentiment was based on the highest votes. Following are the screenshots of confusion matrix and accuracy obtained for dataset-1 and ensemble learning :

### 1. bag of words, 1-gram, classifiers used for voting were : logistic regression, naive bayes, random forest and svm.

```
https://scikit-learn.org/stable/modules/preprocessing.html
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python bag-of-words-evc.py
tcmalloc: large alloc 1216774144 bytes == 0x5627a05dc000 @ 0x7fe10cb76001 0x7fe1096fa54f 0x7fe1
8 0x7fe10974eb17 0x7fe1097ed203 0x56279cc51544 0x56279cc51240 0x56279ccc5627 0x56279cc52afa 0x56
d00 0x56279ccbf9ee 0x56279cc52bda 0x56279ccc1737 0x56279ccbf9ee 0x56279cc52bda 0x56279ccc4d00 0x
bf9ee 0x56279ccbf6f3 0x56279cd894c2 0x56279cd8983d 0x56279cd896e6 0x56279cd61163 0x56279cd60e0c
b95ebf7 0x56279cd60cea
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function
tue_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Ple
get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: ConvergenceWarning
failed to converge (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
  https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
  https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
[[1067  355]
 [ 398  980]]
0.7310714285714286
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1#
```

### 2. positive-negative frequency, 4-grams, classifiers used for voting were : logistic regression, naive bayes and random forest.

```
ols Help All changes saved
+ Code + Text
Terminal X pos-neg-evc.py
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-1# python pos-neg-evc.p
/usr/local/lib/python3.7/dist-packages/sklearn/preprocessing/_label.py:98: DataConversion
column-vector y was passed when a 1d array was expected. Please change the shape of y to (
), for example using ravel().
y = column_or_1d(y, warn=True)
/usr/local/lib/python3.7/dist-packages/sklearn/preprocessing/_label.py:133: DataConversion
column-vector y was passed when a 1d array was expected. Please change the shape of y to
()), for example using ravel().
y = column_or_1d(y, warn=True)
[[1390  621]
 [ 495 1494]]
0.721
```

Following are the screenshots of confusion matrix and accuracy obtained for dataset-2 and positive-negative frequency :

1. positive-negative frequency, 1-gram, logistic regression.

Run: pos-neg-freq

```
[ 1 -1 -1 ... -1 -1 -1]
[[1774   66]
 [ 267  202]]
0.855781723689909
```

2. positive-negative frequency, 2-grams, logistic regression.

Run: pos-neg-freq

```
[-1 -1 -1 ... -1
[[1756   66]
 [ 273  220]]
0.855781723689909
```

Process finished

3. positive-negative frequency, 3-grams, logistic regression.

Run: pos-neg-freq

```
[-1 -1 -1 ...  1 -1 -1]
[[1744   77]
 [ 257  231]]
0.855348635773062
```

Process finished with ex

4. positive-negative frequency, 4-grams, logistic regression.

Run: pos-neg-freq

```
[-1 -1 -1 ... -1 -1 -1]
[[1788   52]
 [ 253  216]]
0.8679081853616284
```

Process finished with ex

5. positive-negative frequency, 3-grams, naive bayes.

Run: y = column\_or\_1d(y,

```
[[1640  181]
 [ 137  351]]
0.8622780424426159
```

Process finished with

6. positive-negative frequency, 4-grams, naive bayes.

```
Run: pos-neg-freq-nb x
y = column_or_1d(y, warn=False)
[[1701 139]
 [ 131 338]]
0.8830662624512776
Process finished with exit code 0
```

7. positive-negative frequency, 3-grams, random forest.

```
Run: pos-neg-freq-nb x
y = column_or_1d(y, warn=False)
[-1 -1 -1 ... 1 1]
[[1702 119]
 [ 178 310]]
0.8713728886964054
Process finished with exit code 0
```

8. positive-negative frequency, 4-grams, random forest.

```
Run: pos-neg-freq-nb x
y = column_or_1d(y, warn=False)
[-1 -1 -1 ... 1 -1 -1]
[[1754 86]
 [ 179 290]]
0.8852317020355132
Process finished with exit code 0
```

9. positive-negative frequency, 3-grams, SVM.

```
Run: pos-neg-freq-nb x
y = column_or_1d(y, warn=False)
[[1708 113]
 [ 190 298]]
0.8687743611953227
Process finished with exit code 0
```

10. positive-negative frequency, 4-grams, SVM.

```
Run: pos-neg-freq-nb x
y = column_or_1d(y, warn=False)
[[1767 73]
 [ 192 277]]
0.8852317020355132
Process finished with exit code 0
```

Following are the screenshots of confusion matrix and accuracy obtained for dataset-2 and bag of words :

1. bag of words, 1-gram, logistic regression.

```

    ret = eng_inst.evaluate()
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# python bag-of-words-lr.py us
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_fe
ture_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please us
get_feature_names_out instead.bda 0x55c524b61737 0x55c524b5f9ee 0x55c524af2bda 0x55c524b64d00 0x55c52
warnings.warn(msg, category=FutureWarning)4c2983d 0x55c524c296e6 0x55c524c01163 0x55c524c00e0c 0x7f5
[[1782   58]5c524c00ceas, local_dict=scope)
[[ 133   336]]iment-analysis-project/sentiment-analysis-of-dataset-2/code# e 824, in evalua/content/ser
0.9172802078822001-analysis-project/sentiment-analysis-of-dataset-2/code/content/sentim
-----
```

## 2. bag of words, 1-gram, multinomial naive bayes.

Terminal X bag-of-words-nb.py X

```

/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# python bag-of-words-nb.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_
ture_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please
get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
[[1775   65]
 [ 154   315]]
0.9051537462104807
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code#
```

## 3. bag of words, 1-gram, random forest.

Terminal X bag-of-words-rf.py X

```

py
py
py
...
09 GB available
[0] 0:python3*          "63caca66fa19" 13
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# python bag-of-wo
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Fun
ture names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.
get_feature_names_out instead.utureWarning)ta.csv', encoding='ISO-8859-1',na_filter=True,
warnings.warn(msg, category=FutureWarning)d
[[1840   0]]14682-analysis-project/sentiment-analysis-of-dataset-2/code/content/sentim
[ 469   0]]187007-analysis-project/sentiment-analysis-of-dataset-2/code/content/sentim
0.7968817669987007-analysis-project/sentiment-analysis-of-dataset-2/code/content/sentim
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code#
```

## 4. bag of words, 1-gram, SVM.

Terminal X bag-of-words-svm.py X

```

/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# python bag-of-words-svm.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_
ture_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please
get_feature_names_out instead.ine 37, in <module>
warnings.warn(msg, category=FutureWarning)ta.csv', encoding='ISO-8859-1',na_filter=True,na_values=
[[1804   36]s=["1": literal_eval])ified 2/data.csv', encoding='ISO-8859-1',na_filter=True,na_values=
[ 196   273]]me 'literal_eval' is not defined
0.8995236032914682-analysis-project/sentiment-analysis-of-dataset-2/code/content/sentim
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code#
[0] 0:bash*          "63caca66fa19" 13:27 23-Nov
```

Following are the screenshots of confusion matrix and accuracy obtained for dataset-2 and TF-IDF :

## 1. TF-IDF, 1-gram, logistic regression.

Terminal X tf-idf-lr.py X

```

/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# pyt
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: Futur
_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will
ease use get_feature_names_out instead.
warnings.warn(msg, category=FutureWarning)
[[1817   23]
 [ 216   253]]
0.8964919878735383
/content/sentiment-analvsis-project/sentiment-analvsis-of-dataset-2/code#
```

## 2. TF-IDF, 1-gram, random forest.

```
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# python tf-idf-rf.py
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function
_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2
  ease use get_feature_names_out instead.
  warnings.warn(msg, category=FutureWarning)
[[1840      0]timent-analysis-project/sentiment-analysis-of-dataset-2/code# python tf-idf-nb.py
 [ 469      0]iment-analysis-project/sentiment-analysis-of-dataset-2/code# tureWarning: Fun/cont
0.7968817669987007
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# █
```

### 3. TF-IDF, 1-gram, SVM.

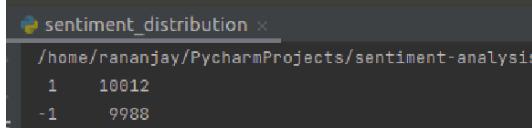
```
warnings.warn(msg, category=FutureWarning)t-analysis-of-dataset-2/code# python tf-idf-nb.py
[[1812     28]]iment-analysis-project/sentiment-analysis-of-dataset-2/code# tureWarning: Fun/cont
 [ 204   265]87007
0.8995236032914682-analysiss-project/sentiment-analysis-of-dataset-2/code/content/sentim
/content/sentiment-analysis-project/sentiment-analysis-of-dataset-2/code# █
[0] 0: bash* "63cac66fa19" 14:58 23-No
```

---

## Data Collection

- Source of dataset-1 : [link](#).
- Source of dataset-2 : [link](#).

Following is the screenshot of distribution of positive(represented as 1) and negative sentiment tweets(represented as -1) in the dataset-1 :



```
sentiment_distribution x
/home/rananjay/PycharmProjects/sentiment-analysis
 1    10012
 -1   9988
```

We can see that both positive and negative sentiment tweets are around 50% of the whole dataset. It's expected since tweets are not from any specific domain.

Following is the screenshot of distribution of positive(represented as 1) and negative sentiment tweets(represented as -1) in the dataset-2:



Percentage of negative and positive sentiment tweets is 79.525 and 20.475, respectively. Since tweets are related to airlines it is expected that most of the tweets are from customers facing some issues and thus there are more negative sentiment tweets.

Following is the screenshot of distribution of positive(represented as 1), negative(represented as -1) and neutral sentiment(represented as 0) tweets in the airline dataset (the complete dataset-2 which doesn't exclude neutral sentiments):



```
-1    9178
 0    3099
 1    2363
```

Percentage of negative, neutral and positive tweets are 62.69, 21.16 and 16.14, respectively.

## Evaluation of Solution

Formula for some of the metrics which will be used to evaluate the results :

$$\begin{aligned} \text{Accuracy} &= \frac{\text{correct predictions count}}{\text{total number of data}} \\ \text{Precision} &= \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \\ \text{Recall} &= \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \\ F_1 \text{ score} &= \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}} \end{aligned}$$

Note that  $F_1$  and  $F$  score refers to the same formula and can be used interchangeably.

For the naive solution (runned on dataset-1) maximum accuracy of 63.6 percent was obtained using a support vector machine. For a very naive solution it's a good result since one of the papers I went through got accuracy around 75 percent after doing all pre-processing, feature extraction. It's interesting to observe however that this naive solution works well in terms of false negatives. In the worst case it's 14.55 percent and in the best case it's 6 percent. So there are very few cases where it says that a tweet has negative sentiment but it's positive in reality.

After I implemented all the techniques that I thought on paper, for dataset-1 maximum accuracy of 0.751 was obtained with bag of words (2-grams, logistic regression) and for dataset-2 maximum accuracy of 0.917 was obtained with bag of words (1-gram, logistic regression). It's expected that accuracy for dataset-2 would be higher because all of it's tweets are related to airlines and thus training would produce better classifiers compared to dataset-1 where tweets don't belong to any specific domain.

For the preprocessing part, following example was runned to illustrate here :

tweet="RT @dynamo #amazing wowwww Thanks!!! for the trip, It was so gooood 🤘"

After preprocessing, the list of word obtained was :

```
['amaze', 'wow', 'thanks', 'thanks', 'trip', 'good', 'thumb', '_amaze_wow', '_wow_thanks', '_thanks_thanks', '_thanks_trip', '_trip_good', '_good_thumb', '_amaze_wow_thanks', '_wow_thanks_thanks', '_thanks_thanks_trip', '_thanks_trip_good', '_trip_good_thumb', '_amaze_wow_thanks_thanks', '_wow_thanks_thanks_trip', '_thanks_thanks_trip_good', '_thanks_trip_good_thumb']
```

In above example we can see that "#amazing" has been converted to "amaze" after removing hashtag and doing lemmatization, "wowwww" has been spell corrected to "wow", "thanks" is repeated twice since in original tweet it had exclamation mark at the end, n-grams till size 4 has been created (example : '\_wow\_thanks\_thanks\_trip'), emoji has been replaced by its corresponding word.

Note that the sizes of dataset I choose while running any code was maximum possible for resources available on google colab. Dataset size for different code varied from 7000 (tf-idf, 4-gram, logistic regression) to 20000 (positive-negative frequency, 1-gram, logistic regression).

It was observed that accuracy increased with more tuple sizes (n-grams). Let's see some of the results obtained from the same dataset sizes and different n-grams for positive-negative frequency for dataset-1.

Table-1

	Model	Accuracy
1.1	training size=16000, test size=4000, 1-gram, SVM	0.6925
1.2	training size=16000, test size=4000, 2-grams, SVM	0.6942
1.3	training size=16000, test size=4000, 3-grams, SVM	0.6942
1.4	training size=16000, test size=4000, 4-grams, SVM	0.7072

Table-2

	Model	Accuracy
1.1	training size=16000, test size=4000, 1-gram, logistic regression	0.7055
1.2	training size=16000, test size=4000, 4-grams, logistic regression	0.7182

Some of the results obtained from the bag of words model for dataset-1.

Table-3

	Model	Accuracy
1.1	training size=11200, test size=2800, 1-gram, logistic regression	0.7246
1.2	training size=11200, test size=2800, 2-grams, logistic regression	0.7510

Table-4

	Model	Accuracy
1.1	training size=11200, test size=2800, 1-gram, logistic regression	0.7271
1.2	training size=11200, test size=2800, 2-grams, logistic regression	0.7510

It's not correct to compare different n-grams with very different dataset sizes. While running the algorithms, memory limit was exceeded in google colab for 4-grams but it ran for 1-gram for the same large input dataset size.

Some of the results obtained from the positive-negative frequency for dataset-2.

Table-5

	Model	Accuracy
1.1	training size=9232, test size=2300, 3-grams, naive bayes	0.8622
1.2	training size=9232, test size=2309, 4-grams, naive bayes	0.8830

From the above tables we can say that accuracy increased around 2% in some cases when the number of tuples (n-grams) were increased.

Following table represents  $F_1$  scores for the cases where best accuracy was obtained for each dataset :

Table-6

	Model	$F_1$
dataset-1	training size=11200, test size = 2800, bag of words, 2-grams, logistic regression	0.7573
dataset-2	training size=9232, test size=2309, bag of words, 1-gram, logistic regression	0.7786

Following table has results of ensemble learning for bag of words and 1-gram for dataset-1:

Table-7

	Model	Accuracy
1.1	training size=11200, test size = 2800, hard voting (logistic regression, naive bayes, random forest, SVM)	0.7310
1.2	XGBoost	0.6867

Let's compare the result of hard voting with individual algorithms:

Table-8

	Model	Accuracy
1.1	training size=11200, test size = 2800, logistic regression	0.7246
1.2	training size=11200, test size = 2800, naive bayes	0.7121
1.3	training size=11200, test size = 2800, random forest	0.7142
1.4	training size=11200, test size = 2800, SVM	0.7271

We can see that the maximum accuracy in the above table is 0.7271 and by using the hard voting ensemble learning technique it increased to 0.7310.

Also I choose only those algorithms for hard voting whose accuracy are comparable. This is because if we choose lets say 25 models which can vote and suppose 10 of them have accuracy 0.72 and 15 of them have accuracy 0.63, chances are that there won't be any improvement from 0.72.

There is another ensemble learning technique called soft voting which uses probability but since SVM doesn't output probability thus it is not possible to use soft voting with it.

### Comparison with other works

Let's compare the results in this report with those of Samuel, Ali, Rahman, Esawi, Samuel[2]. There is a study of the effect of accuracy on the number of characters in COVID-19 related tweets.

Following are the confusion matrices (in the paper) for Naive Bayes for various tweet lengths :

Number of characters less than 77, accuracy=0.9143

	Negative	Positive
Negative	34	1
Positive	5	30

Number of characters less than 120, accuracy=0.5714

	Negative	Positive
Negative	34	1

Positive	29	6
----------	----	---

Following are the confusion matrices (in the paper) for Logistic Regression for various tweet lengths :

Number of characters less than 77, accuracy=0.7429

	Negative	Positive
Negative	30	5
Positive	13	22

Number of characters less than 120, accuracy=0.52

	Negative	Positive
Negative	21	14
Positive	19	16

Following is the confusion matrix obtained for positive-negative frequency, unigram, multinomial naive bayes for dataset-1 in my work :

Accuracy : 0.7022

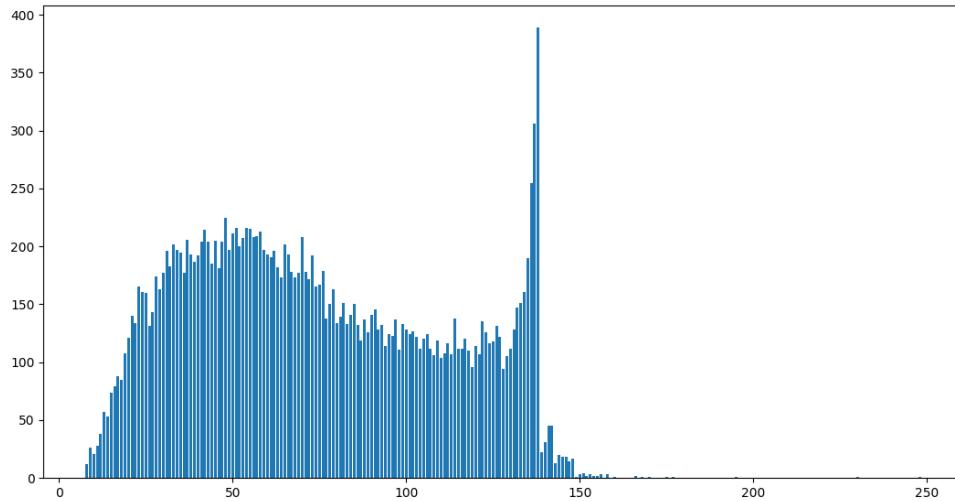
	Negative	Positive
Negative	1357	630
Positive	561	1452

Following is the confusion matrix obtained for bag of words, 2-grams, logistic regression for dataset-1 in my work :

Accuracy : 0.751

	Negative	Positive
Negative	1015	376
Positive	321	1088

Following is the distribution of tweet length in the dataset :



Out of 20000 tweets, 8564 have length greater than 77 and 11434 have length less than 77.

Let's compare the F scores (2nd column has F score for the research paper and 3rd column has F score for my work) :

Naive Bayes	0.9090, 0.2857	0.7091
Logistic Regression	0.7096, 0.4923	0.7573

Note that the 2nd column represents the F score of both : number of characters less than 77 and less than 120 (for example 0.9090, 0.2857 for naive bayes) in the research paper. We can say that Logistic Regression performed much better in my case. Also note that dataset-1 is generic whereas the dataset used by authors is related to COVID-19. A test with a dataset related to some single domain is expected to give better F score since words which play an important role in determining the sentiment would be less and occur a lot more time than in dataset which contains tweets from every domain.

Following table represents accuracy obtained by Kharde, Sonawane[3] and by me for naive bayes :

Algorithm	Kharde, Sonawane	dataset-1	dataset-2

Naive Bayes (unigram)	74.56	71.21	86.53
Naive Bayes (bigram)	76.44	72.54	86.4
Naive Bayes (trigram)	75.41	72	86.22

The results in the column of dataset-1 is for bag of words and results in the column of dataset-2 is for positive-negative frequency.

We can see that there is a slight decrease in accuracy for dataset-2 and dataset-1 with increase in n-grams. Similar patterns can be observed for Kharde, Sonawane when we compare bigram and trigram. There is an increase in the number of tokens with n-grams and when a sufficiently large dataset is not available the tokens formed with higher n-grams would be present in less number of tweets. It can then lead to poor training and those tokens would just skew the true results while testing.

Following table represents accuracy obtained by Kharde, Sonawane and by me for SVM :

Algorithm	Kharde, Sonawane	dataset-1	dataset-2
SVM with unigram	76.68	72.71	85.57
SVM with bigram	77.73	73.69	85.57

The results in the column of dataset-1 is for bag of words and results in the column of dataset-2 is for positive-negative frequency.

Similar to Kharde, Sonawane we can see that for dataset-1 we obtained more accuracy for bigram than unigram.

From the results it seems that type of algorithm plays an important role in determining whether accuracy will increase with increase in n-grams. It is observed from the screenshot section (“Code and Output”) that logistic regression and support vector machines usually give more increase in accuracy with increase in n-grams than naive bayes. For example, for bag of words, 72.46 accuracy was obtained for unigram and 75.10 for bigram with logistic regression whereas for naive bayes 71.21 accuracy was obtained for unigram and 72.54 for bigram.

Following table compares results obtained from 3 different sentiment analysis tools[6] :

	Proposed System	Alchemy API	SVM (Pre-processed tweets)	SVM, dataset-1
Accuracy	59.85%	58.87%	64.95%	73.69%
Precision	53.65%	40.82%	66.54%	71.33%
F-measure	0.48	0.43	0.57	70.21%

Column ‘SVM, dataset-1’ is for preprocessed tweets with bag of words and bigram. We can clearly see that I obtained better results compared to those in the first three columns. Note that dataset-1 contains almost 50% percent positive tweets and 50% negative tweets (to be precise : 10012 positive tweets and 9988 negative tweets from randomly chosen 20000 tweets from dataset-1).

For the airline dataset, I also ran algorithms for tweets with three categories : positive, neutral and negative sentiment. Results can be compared with those obtained by Taufic, Touhid[7] who also experimented with the same dataset.

Following is a comparison for the bag of words.

Algorithm	Accuracy Taufic, Touhid	Accuracy obtained by me (unigram, bigram)
SVM	0.77	0.78, 0.735
Multinomial Naive Bayes	0.74	0.762, 0.765
Logistic Regression	0.77	0.783, 0.78

Following is a comparison for the TF-IDF.

Algorithm	Accuracy Taufic, Touhid	Accuracy obtained by me (unigram, bigram)
SVM	0.77	0.743, 0.716
Multinomial Naive Bayes	0.70	0.686, 0.676
Logistic Regression	0.77	0.755, 0.766

We can see that for TF-IDF, the results of Taufic, Touhid are slightly better than my results. I have obtained slightly better results for SVM, Multinomial Naive Bayes, Logistic Regression in bag of words. Note that for TF-IDF with bigram I am getting less accuracy because I used less dataset size due to the memory limit of google colab (for example I used 9500 tweets for training and testing in TF-IDF with bigram and SVM whereas total dataset size is 14640).

Since I obtained better result for bag of words than the authors, let's compare weighted F scores (author has also used weighted F scores since dataset is imbalanced) for the same :

	F score Taufic, Touhid	F score obtained by me (unigram)
SVM	0.75	0.769
Multinomial Naive Bayes	0.72	0.745

Logistic Regression	0.77	0.780
---------------------	------	-------

We can clearly see that I obtained better F scores than the authors for bag of words.

#### Comparison with python sentiment analysis tools

	Accuracy from Vader	Accuracy from Textblob	Accuracy obtained by me
dataset-1	0.52405	0.6192	0.751 (bag of words, logistic regression, bigram)
dataset-2 (with positive,negative and neutral tweets)	0.490	0.464	0.78 (bag of words, logistic regression, unigram)

Following are confusion matrices obtained for vader (-1 refers to negative sentiment, 0 refers to neutral, 1 refers to positive sentiment):

dataset-1

	-1	0	1
-1	4243	2548	3197
0	0	0	0
1	1016	2758	6238

dataset-2

	-1	0	1
-1	4018	1445	3715
0	350	1027	1722
1	65	159	2139

Following are confusion matrices obtained for textblob :

dataset-1

	-1	0	1

-1	6781	0	3207
0	0	0	0
1	4409	0	5603

dataset-2

	-1	0	1
-1	3231	3287	2660
0	331	1759	1009
1	104	453	1806

Following are confusion matrices obtained from my solution where I divided data into 80-20 for training and testing :

dataset-1

	-1	0	1
-1	1015	0	376
0	0	0	0
1	321	0	1088

dataset-2

	-1	0	1
-1	1629	156	53
0	192	332	69
1	94	70	333

From above we can see that I have obtained the highest accuracy in all cases. If we compare F score for highest accuracy case : F score for Textblob is 0.559 for dataset-1 and we have obtained F score of 0.757.

F score, dataset-1

Textblob	My code
0.59	0.757

### **other techniques used for improving performance**

I have used the effect of exclamation marks and utilized information of words with underscore and hashtags which can increase the accuracy. A large number of tweets related to politics and sports have hashtags, words with underscores and usually they contain words other than stop words and keeping them can be useful compared to the cases where we just remove any token having character other than alphabet. Also, every technique is not used in every research paper. For example some researchers did not use n-grams, however I have tried to use each and every preprocessing technique which can improve the results.

## Conclusion

- Based on comparison using different parameters (accuracy, precision,  $F_1$  score, tweets with different numbers of characters) with 4 different research papers, I obtained better results than 2 of them and quite close to the other two. Note that dataset-1 has tweets from every domain and almost equal number of positive and negative tweets are present. If any dataset has an unequal number of positive-negative tweets or is from some specific domain (say COVID-19, movie review) then it's likely to give better results for some parameters as we observed with airline dataset (dataset-2). I obtained 0.917 accuracy for dataset-2 and 0.751 accuracy for dataset-1. However the  $F_1$  scores, 0.75 for dataset-1 and 0.77 for dataset-2 are very close.
- I obtained better results in terms of F score and accuracy if we compare numbers obtained with vader and textblob.
- $F_1$  score of 0.7614 was obtained for dataset-1 and  $F_1$  score of 0.8045 was obtained for dataset-2.
- Maximum accuracy of 0.751 was obtained for dataset-1 and maximum accuracy of 0.917 for dataset-2.
- Accuracy for dataset-1 is less because it's tweets are not from any specific domain. Dataset-2 contains airline dataset and so some words would occur more frequently and thus we got more accuracy for it.
- Bag of Words and TF-IDF gave comparable results and performed significantly better than Positive-Negative frequency. This is because in Bag of Words and TF-IDF we are considering all words whereas in Positive-Negative frequency all information is put into 2-dimension. Bag of words and TF-IDF have more time and memory complexity compared to positive-negative frequency.
- Joining multiple words (n-grams) gave better results compared to unigrams (1-gram).

## **Future Work**

Currently I have thought about following ways to make further improvements :

- We can replace all synonym words with a common word in each tweet.
- We used n-grams for taking account of negation (for example : “not sad”). However we can try using antonyms (i.e replace “not sad” with “happy”) along with n-grams (for cases like “not very happy”, “not” and “happy” have one word between them and thus we can either use antonym for words few distance apart or we can use higher n-grams like 3-grams, 4-grams).
- Trying to improve preprocessing and better ways of representing the data mathematically.

## References

- [1] Alonso, Vilares, Gómez-Rodríguez, Vilares, [Sentiment Analysis for Fake News Detection](#), 2021
- [2] Samuel, Ali, Rahman, Esawi, Samuel, [COVID-19 Public Sentiment Insights and Machine Learning for Tweets Classification](#), 2020
- [3] Kharde, Sonawane, [Sentiment Analysis of Twitter Data: A Survey of Techniques](#), 2016
- [4] Alsaeedi, Khan, [A Study on Sentiment Analysis Techniques of Twitter Data](#), 2019
- [5] [Twitter character limit](#)
- [6] Chong, Selvaretnam, Soon, [Natural Language Processing for Sentiment Analysis](#), 2014
- [7] Taufic, Touhid, [A Comparative Study of Sentiment Analysis Using NLP and Different Machine Learning Techniques on US Airline Twitter Data](#), 2021