Aarifah Ullah

Prof. Kasia Hitczenko

CSCI 4907.82

3 February 2024

<div align="center">CSCI 4907.82 Natural Language Understanding Homework 1</div>

**Part 1: Preparing Data**

1. What are the dimensions of your feature matrix?
    a. The matrix is 5000 X 9084, meaning there are 5000 tweets and 9084 unique words that appear within those tweets. This makes sense as the csv file is after all named "Tweets_5K.csv"
2. What is the value of X[1460][1460]?
    a. The value is 0.
3. What does this value mean? What feature does the 1460th column represent?
    a. This value is the count of the 1460th word in the matrix in the 1460th document of the matrix. So, the 1460th word is "bristol" and we are asking how many times it appears in document 1460. This is 0 because that document says, "Happy Mother's Day!"

**Part 2: Implementing Naive Bayes**

1. Calculate whether the Naive Bayes classifier you trained would classify the following tweet as positive, neutral, or negative: "Happy birthday Annemarie". You should do this by hand and show your work, but you can use code to get the relevant probabilities (basically, the goal of this question is to make sure that you understand how Naive Bayes works)
    a.
2. Report the model's accuracy on the unseen test set. How does this compare to a classifier that always outputs the most frequent category in the training set (what is that classifier's accuracy)?
    a. The naive bayes with just the preprocessing steps had an F1 score of .5554. An F1 score closer to 1 indicates a very good classifier while closer to zero is a poor classifier. The Naive Bayes classifier seems to do okay-ish.

       The accuracy is .54575 or 54.575%, meaning the classifier gets the right answer about half the time. Compared to a classifier that just always outputs the most frequent category, this Naive Bayes classifier has a lower accuracy. This is because the other classifier just answers the most frequent class and is right about most of the time, but it is a poor classifier since it does not know much about the data. (It's precision and recall are perhaps very low).

**Part 3: Implementing Logistic Regression**

1. In class, we learned how to do binary logistic regression between 2 options. Here, there are three possible classifications, so sklearn uses a "one vs. rest" scheme where it learns a

binary logistic regression model for each possible label (i.e., one logistic regression which learns to separate positive from non-positive tweets, a second that learns to separate negative from non-negative tweets, and a final logistic regression that learns to separate neutral from non-neutral tweets). How many parameters did your multiclass logistic regression model learn?

    a. Parameters and features are different. Parameters refers to the document's attributes that set it apart and help to place that document into groups known as features. Because we implemented a bag of words representation, the words themselves are what defines the documents, as in the document can be described by all the words it has. So, in a binary logistic model using the bag of words representation, it learns the whole list of words as the parameters. Because the multinomial logistic regression model employed binary logistic regression model three times (positive v. rest, negative v. rest, neutral v. rest), it still learned the same number of parameters, but repeated its steps three times.

2. Report the model's accuracy on the unseen test set. How does this compare to a classifier that always outputs the most frequent category in the training set and the Naive Bayes classifier from Part 2?

    a. The logistic regression model has an F1 score of 55.54% and an accuracy of .58 or 58%, meaning that the model correctly assesses the class 58% of the time. This is a little better than the naive bayes classifier that had an accuracy rate of 54.575% using just basic pre-processing steps and the bag of words model. Compared toa classifier that outputs the most frequent category, the logistic regression model is less accurate for the same reasons mentioned when describing the naive bayes classifier. The other classifier is not very good but it has a high accuracy rate when it guesses the most frequent class all the time.
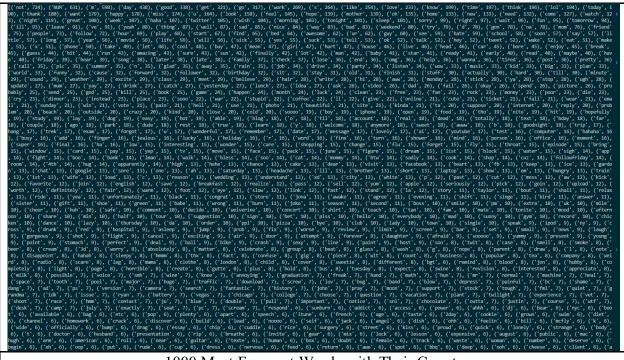
## Part 4: Implementing More Elaborate Pre-Processing

1. How does this impact the performance of your Naive Bayes and Logistic Regression classification results?

    a. More pre-processing improves results. After implementing each processing step, I re-ran both models and observed how both the accuracy and F1 score improved. The final accuracy and logistic numbers are

## Part 5: Your turn!

1. Propose and implement at least one addition to the workflow that you think will improve performance (but please still use Naive Bayes or Logistic Regression). Specifically, you can add a pre-processing step, add a feature (or class of features), etc. Explain why you chose this feature and why you think it might help performance. Report the classifier's accuracy on the test set. Did this help as expected? Note: You will not be graded on whether your change actually improved performance.

    a. My suggestion to improve performance is to make a separate list of stop words that add onto Spacy's standard list of stop words. This would help eliminate words more related to tweets that are perhaps taking up space in the 1000 most frequent words lists and are not giving us clues on the actual sentiment.

b. After generating a list of the 1000 most frequent words that occur after more elaborate pre-processing including removal of stop words and lemmatization, it is clear to see that some "words" remain.

c. For example:

  i. "ve" is not a word but most likely a conjunction that got separated from the root word with the removal of punctuation "ve" is one of the most frequent words, in the top of the list with a count of 91.

  ii. "s" ~ 264, "m" ~ 598 most likely from "I'm," "d" ~78

d. I made a custom list of common leftover words I saw and include that along with my stop word function on line 103.I also used CountVectorizer()'s own stopword feature to handle any missed word .isstop did not handle.



1000 Most Frequent Words with Their Counts

2. Finally, take the best performing model and look at the tweets that were incorrectly categorized. Do you observe any patterns in what the model is making mistakes on? What do you think could be done to further improve results? Provide a sample of tweets as well as their true label and their predicted label to justify your answers. Note: You do not need to know how to implement the proposed change.

  a. Logistic regression performed the best with an estimated accuracy of ~60%. It seems to wrongly classify texts that are difficult to tell as a person. It would depend on the context. It also got wrong texts with words and phrases that are used in other classes as well. For example "soooooo..." can be used in a positive sense, but in line 1, the logistic regression model didn't realize it was meant in a negative sense. I would suggest working with more data to improve the results.

```
import pandas as pd #read/writing csv file
        textID                                               text sentiment
0     cb774db0d1            I`d have responded, if I were going   neutral
1     549e992a42         Sooo SAD I will miss you here in San Diego!!!  negative
2     088c60f138                          my boss is bullying me...  negative
3     9642c003ef               what interview! leave me alone  negative
7     50e14c0bb8                                 Soooo high   neutral
...          ...                                                ...       ...
3989  7fd2cc86b0                 I just wanna be better already  positive
3990  9704c2d674   : have fun focus grouping! i have no participa...  positive
3991  35b97b5ac2       Thanks! Am trying.. (keeping my chin up, ouch!)   neutral
3995  ed57123f1c  Tinkered with open-source Virtualbox &Win7 yet...  neutral
3999  b59a68c7f5   also: is there a strategic IT plan that maps ...   neutral
```

List of Tweets Logistic Regression Incorrectly Classified

3. Roughly how much time did you spend on this homework?
   a. I have spent about 16 hours completing this homework. I haven't coded in a bit, so it took me some time to set up my computer, including downloading python, etc. I also had some trouble downloading the necessary NLP libraries.

|  | Naive Bayes | | Logisitic Regression | |
|---|---|---|---|---|
|  | Accuracy | F1 | Accuracy | F1 |
| Pre-Processing | 54.575% | 55.54% | 58% | 55.54% |
| More Elaborate Pre-Processing | ~57.35% | 57.38% | 59.625% | 57.83% |
| My Own Implemented Suggestion | 56.85% | 56.93% | 58.65% | 56.93% |

**Consulted Sources**

When implementing the Bag of Words matrix:

- https://spotintelligence.com/2023/05/17/countvectorizer/
- https://spotintelligence.com/2022/12/20/bag-of-words-python/

Splitting Training and Testing sets:

- https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
- https://www.sharpsightlabs.com/blog/scikit-train_test_split/

Naive Bayes Implementation:

- https://scikit-learn.org/stable/modules/naive_bayes.html
- https://www.datacamp.com/tutorial/naive-bayes-scikit-learn
- https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html

Logistic Regression Implementation:

- https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html
- https://towardsdatascience.com/logistic-regression-using-python-sklearn-numpy-mnist-handwriting-recognition-matplotlib-a6b31e2b166a

More Pre-Processing:

- https://www.digitalocean.com/community/tutorials/python-remove-spaces-from-string
- https://blog.enterprisedna.co/python-remove-punctuation-from-string/
- https://www.geeksforgeeks.org/python-pos-tagging-and-lemmatization-using-spacy/
- https://realpython.com/natural-language-processing-spacy-python/
- Basics of CountVectorizer | by Pratyaksh Jain | Towards Data Science