

Springboard Data Science 6 months program
Capstone Project 2

Sentiment analysis of Tweets about the movie *Joker*

Coline Zeballos

2019

Table of Contents

Proposal	3
The problem I want to solve	3
What is a Sentiment Analysis and why is it useful?	3
My potential client and why he cares about this problem	3
The data I will be using and how I will acquire it.....	3
How I will solve this problem	3
My deliverables.....	4
Methodology	5
Section A: Preparing the Test set.....	6
Step A.1: Getting the authentication credentials	6
Step A.2: Authenticating our python script.....	6
Step A.3: Creating the function to build the Test set.....	6
Section B: Preparing the Training set	8
Section C: Pre-processing Tweets in the Data Sets	10
Step C.1: What matters and what doesn't matter in Sentiment Analysis?.....	10
Step C.2: A word about the importance of normalizing/pre-processing	10
Step C.3: Explaining the pre-processing steps and their importance	10
Section D: Naive Bayes Classifier	12
Why Naive Bayes Classifier?	12
Other possible classifier	12
Step D.1: Build a vocabulary/list of words in our training data set.....	12
Step D.2: Match tweet content against our vocabulary	12
Step D.3: Build our word feature vector	12
Step D.4: Training the classifier	13
Section E: Testing the model	14
Section F: Analyzing the performance of the model.....	15
Step F.1: Concatenate pre-processed Tweets and labels	15
Step F.2: Check 10 pre-processed Tweets and their attributed label	16
Step F.3: Model performance	17
Step F.4: Wrap up observations.....	17
Step F.5: Possible improvements	17
Conclusion	18

Proposal

The problem I want to solve

I would like to perform a sentiment analysis on Tweets about the late blockbuster release: Joker, with Joaquin Phoenix.

What is a Sentiment Analysis and why is it useful?

Sentiment Analysis represents the use of Natural Language Processing to determine the attitude, opinions and emotions of a speaker, writer or other subject within an online mention. In other words, it is the process of determining whether a piece of writing is positive or negative. A human is able to recognize and classify a text into positive or negative. However, a computer not but can learn to do so.

My potential client and why he cares about this problem

Out in the USA beginning of October 2019 and throughout theaters in the world later on, no doubt that Joker has divided its audience. Whether people love or hate it, the reactions were numerous. I saw it myself at the movies in Switzerland as soon as it came out, and felt exactly this way: divided. The movie received a Golden Lion win at the Venice Film Festival, I had high expectations. The complex personality of the Joker was one of the reasons I wanted to see the movie. But after visioning it, I realized it had a strong impact on me, more than expected and not only positive. The critics I read afterwards were also two-folds, both positive and negative. Therefore, I thought it would be interesting to train a classification model on Tweets on the topic of Joker.

Different types of people could be interested to know the proportion of positive vs negative tweets on this topic: fans of Joaquin Phoenix and of The Dark Knight Rises, owners of movies, producers, script writers, psychologists and psychiatrists.

The data I will be using and how I will acquire it

I will be using the Twitter API to collect a Test set based on keywords. A function will return a list of tweets that contain our keywords selected. Each tweet's text will see itself attributed a label ('positive' or 'negative') to classify each tweet as positive or negative. The Training set will be downloaded because it has to be labelled into 'positive' or 'negative' on a big amount of tweets. The Training set is critical to the success of the model since our model will "learn" how to do create a sentiment analysis based on the Training set.

How I will solve this problem

The steps that we will follow to perform the sentiment analysis are:

- acquiring a Test set
- acquiring and preparing a Training set
- pre-processing tweets in both Test set and Training set
- building a vocabulary/list of words in our training data set
- matching tweet content against our vocabulary
- building our word feature vector
- training the classifier
- testing the model
- evaluating the model performance

My deliverables

A Jupyter notebook containing the code, a slide deck , and a final report in PDF.

All deliverables are available on the following Github repository:

<https://github.com/Crosita/Capstone-Project-2>

Methodology

I will be using the Twitter API to collect a Test set based on keywords. A function will return a list of tweets that contain our keywords selected. Each tweet's text will see itself attributed a label ('positive' or 'negative') to classify each tweet as positive or negative. The Training set will be downloaded because it has to be labelled into 'positive' or 'negative' on a big amount of tweets. The Training set is critical to the success of the model since our model will "learn" how to do create a sentiment analysis based on the Training set.

The steps that we will follow to perform the sentiment analysis are:

- Section A: Preparing the Test set
 - Step A.1: Getting the authentication credentials
 - Step A.2: Authenticating our python script
 - Step A.3: Creating the function to build the Test set
- Section B: Preparing the Training set
- Section C: Pre-processing Tweets in the Data Sets (both Test and Training)
- Section D: Naive Bayes Classifier
 - Step D.1: Build a vocabulary/list of words in our training data set
 - Step D.2: Match tweet content against our vocabulary
 - Step D.3: Build our word feature vector
 - Step D.4: Training the classifier
- Section E: Testing the model
- Section F: Analyzing the performance of the model
 - Step F.1: Concatenate pre-processed Tweets and labels
 - Step F.2: Check 10 pre-processed Tweets and their attributed label
 - Step F.3: Model performance
 - Step F.4: Wrap up observations
 - Step F.5: Possible improvements

Section A: Preparing the Test set

Step A.1: Getting the authentication credentials

To do so, I visited the Twitter Developer website, applied for an API access and waited for Twitter to accept. After getting the approval email, I went back to the website and clicked on “Create an app”. I filled in the information asked and finally got Consumer API keys, as well as Access token & access token secret.

Step A.2: Authenticating our python script

I then imported the twitter library, and created a twitter.Api object containing the credentials mentioned previously.

Step A.3: Creating the function to build the Test set

Now we can start on making a function that downloads the Test set that we talked about. Basically, this is going to be a function that takes a search keyword (i.e. string) as an input, searches for tweets that include this keyword and returns them as twitter.Status objects that we can iterate through.

The caveat here, though, is that Twitter limits the number of requests you can make through the API for security purposes. This limit is 180 requests per 15-minute window. This means, we can only get up to 180 tweets using our search function every 15 minutes, which should not be a problem, as our Training set is not going to be that large anyway. For the sake of simplicity, we will limit the search to 100 tweets for now, not exceeding the allowed number of requests. Our function for searching for the tweets (i.e. Test set) will be:

```
def buildTestSet(search_keyword):  
    try:  
        tweets_fetched = api.GetSearch(search_keyword, count = 100)  
  
        print("Fetched " + str(len(tweets_fetched)) + " tweets for the term " + search_keyword)  
  
        return [{"text":status.text, "label":None} for status in tweets_fetched]  
    except:  
        print("Unfortunately, something went wrong..")  
        return None
```

As expected, this function will return a list of tweets that contain our search keyword/s.

Note that we coupled — into a JSON object — every tweet’s text with a label that is NULL for now. This is merely because we are going to classify each tweet as Positive or Negative later on, in order to determine whether the sentiment on the search term is positive or negative, based on the majority count. This is how Sentiment Analysis pragmatically works.

As described in the project introduction, we will be looking for 100 Tweets talking about the movie "Joker", we will therefore search for the keywords "joker" and "movie". The Tweets will be stored in a variable called testDataSet.

```
print(testDataSet[0:4])
```

```
[{'text': 'Finally saw the joker.\n\nWhile the technical\naspects of the movie are stunning and the acting incredib\nle.\n\nThat migh... https://t.co/QFIi10T0oH', 'label': None},\n {'text': '"Joker" became the first R-rated movie to hi\n t $1 billion at the box-office, and there are\u200b rumor\n s of a sequel https://t.co/k6U8yuxsNC', 'label': None},\n {'text': "Would you like to see a sequel to Joaquin Phoeni\n x's #Joker? https://t.co/InEGdlEvME", 'label': None},\n {'text': 'RT @lwool7: They had done this before for other mo\n vies such as:\nBTS Bring The Soul: The Movie\nWeathering\nWith You\nJoker\nSpiderman: Far From...', 'label': None}]
```

These are five tweets that contain our search keywords. Now everything is set. We have our Test set and we can move on to building our Training set.

Section B: Preparing the Training set

In this section, we will also be using our Twitter API instance from the last section. However, we need to get some things out the way first. We will be using a downloadable Training set. The tweets of which were all labeled as positive or negative, depending on the content. This exactly what a Training set is for.

A Training set is critical to the success of the model. Data is which needs to be labeled properly with no inconsistencies or incompleteness, as training will rely heavily on the accuracy of such data and the manner of acquisition.

For this task, we will be using Niek Sanders' Corpus of over 5000 hand-classified tweets, which makes it quite reliable. There's also a catch here. Twitter does not allow storing tweets on a personal device, even though all such data is publicly available. Therefore, the corpus includes a keyword (topic of the tweet), a label and a tweet ID number for every tweet (i.e. row in our CSV corpus). You can get the file containing the corpus through this link of a personal repository: <https://github.com/karanluthra/twitter-sentiment-training/blob/master/corpus.csv>

Let's explore the Corpus dataset:

```
tr_set.head()
```

	apple	positive	126415614616154112
0	apple	positive	126404574230740992
1	apple	positive	126402758403305474
2	apple	positive	126397179614068736
3	apple	positive	126395626979196928
4	apple	positive	126394830791254016

Characteristics of the Corpus:

- 5512 lines
- 3 columns
- Column Topic contains 4 topics

```
: tr_set[tr_set.columns[0]].unique()  
  
: array(['apple', 'google', 'microsoft', 'twitter'], dtype=  
object)
```

- Column Label contains 4 labels

```
tr_set[tr_set.columns[1]].unique()  
  
array(['positive', 'negative', 'neutral', 'irrelevant'],  
dtype=object)
```

The original dataset "tr_set" doesn't contain the Tweets text, but we will fetch it using the API below.

I will be using the API to get the actual tweet text through each tweet's ID number included in the Corpus we have. Since Twitter API has a limit mentioned above, I will make the code rest for 15 minute to download 5000 tweets.

The function to do so does the following:

- two inputs: "corpusFile" – a string path to the Niek Sanders' corpus file I downloaded (containing tweet's **topic**, **label** and **id**), and "tweetDataFile" – a string path to the file we would like to save the full tweets in (containing tweet's **text**, **topic**, **label** and **id**)
- started with empty list *corpus*
- then opened file "corpusFile" and appended every tweet from the file to the list *corpus*
- get the text of tweets based on the ID
- loop through the tweets in *corpus* calling the API on every tweet to get the Tweet.Status object of the particular tweet.
- use status object to get the text associated with it and push it into the "trainingDataSet"
- sleep for a couple of minutes to avoid the API limit mentioned

The tweets retrieved through the API are saved in a pickle file "trainingData".

```
In [13]: trainingData
```

```
Out[13]: [{'tweet_id': '126415614616154112',  
          'label': 'positive',  
          'topic': 'apple',  
          'text': 'Now all @Apple has to do is get swype on the i  
phone and it will be crack. Iphone that is'},  
          {'tweet_id': '126402758403305474',  
          'label': 'positive',  
          'topic': 'apple',  
          'text': "Hilarious @youtube video - guy does a duet wit  
h @apple 's Siri. Pretty much sums up the love affair! ht  
tp://t.co/8ExbnQjY"},  
          {'tweet_id': '126397179614068736',  
          'label': 'positive',  
          'topic': 'apple',  
          'text': '@RIM you made it too easy for me to switch to  
@Apple iPhone. See ya!'},
```

Section C: Pre-processing Tweets in the Data Sets

Before we move on to the actual classification section, there is some cleaning up to do. As a matter of fact, this step is critical and usually takes a long time when building Machine Learning models. However, this will not be a problem in our task, as the data we have is relatively consistent. In other words, we know exactly what we need from it. I will express on this matter later on.

Step C.1: What matters and what doesn't matter in Sentiment Analysis?

Words are the most important part (to an extent that we will talk about in the upcoming section). However, when it comes to things like punctuation, you cannot get the sentiment from punctuation. Therefore, punctuation does not matter to Sentiment Analysis. Moreover, tweet components like images, videos, URLs, usernames, emojis, etc. do not contribute to the polarity (whether it is positive or negative) of the tweet. However, this is only true for this application. For instance, in another application, we could have a Deep Learning image classifier that learns and predicts whether this image that the tweet contains stands for something positive (e.g. a rainbow) or negative (e.g. a tank). When it comes to the technicality, both Sentiment Analysis and Deep Learning fall under Machine Learning. In fact, we can perform Sentiment Analysis through Deep Learning, but that's a story for another day.

Step C.2: A word about the importance of normalizing/pre-processing

Normalization in the NLP context is the process of converting a list of words to a more uniform sequence. By transforming the words into a standard format, later operations can be done on the data without compromising the process. Many pre-processing steps can be taken including: lowercasing, stemming (example: troubled, troubles go into 'troubl'), lemmatization (example: troubled, troubles to into 'trouble'), stop word removal, normalizing, noise removal...

Step C.3: Explaining the pre-processing steps and their importance

1. We start off by our imported libraries. `re` is Python's Regular Expressions (RegEx) library, which takes care of parsing strings and modifying them in an efficient way without having to explicitly iterate through the characters comprising the particular string. We also imported `ntlk`, is the Natural Processing Toolkit, which is one of the most commonly used Python libraries out there. It takes care of any processing that we need to perform on text to change its form or extract certain components from it. The class constructor removes stop words. This is a relatively big topic that you can read up on later, as it is more into Natural Language Processing and less related to our topic.
2. The `processTweets` function loops through all the tweets input into it, calling its neighboring function `processTweet` on every tweet in the list. The latter does the actual pre-processing steps:
 - Removing stop words in English: stop words are a set of commonly used words in a language (a, the, is, are,...). The idea behind removing these words is to remove low information words from the text and focus on the important words instead.
 - Making all the text in lower-case letters. This is merely because, in almost all programming languages, "cAr" is not interpreted the same way as "car". Therefore, it is better to normalize all characters to be lower-case across all our data.
 - URLs and usernames are removed from the tweet. This is for the reasons we disclosed earlier.

- The number sign (i.e. #) is removed from every hashtag, in order to avoid hashtags being processed differently than regular words.
- Duplicate characters are rid off of, in order to ensure that no important word goes unprocessed even if it is spelled out in an unusual way (e.g. “caaaaar” becomes “car”).
- Finally, the tweet’s text is broken into words (tokenized) in order to ease its processing in the upcoming stages. This step is also called text segmentation or lexical analysis.

We apply the pre-processing steps to both the Training Set and the Test set.

Note: Remember, the Tweets' labels in the Training Dataset take 4 values: positive, negative, irrelevant and neutral. We decide at this point to select on Tweets with Label 'positive' and 'negative'.

```
In [18]: print(preprocessedTrainingSet[0:4])

[(['get', 'swype', 'iphone', 'crack', 'iphone'], 'positive'), ([ 'hilarious', 'video', 'guy', 'duet', "'s", 'siri', 'pretty', 'much', 'sums', 'love', 'affair'], 'positive'), ([ 'made', 'easy', 'switch', 'iphone', 'see', 'ya'], 'positive'), ([ '16', 'strangest', 'things', 'siri', 'said', 'far', 'sooo', 'glad', 'gave', 'siri', 'sense', 'humor', 'via'], 'positive')]
```

```
In [36]: print(preprocessedTestSet[0:4])

[(['finally', 'saw', 'joker', 'technical', 'aspects', 'movie', 'stunning', 'acting', 'incredible', 'migh...'], None), ([ '...', 'joker', '...', 'became', 'first', 'r-rated', 'movie', 'hit', '1', 'billion', 'box-office', 'are\u200b', 'rumors', 'sequel'], None), ([ 'would', 'like', 'see', 'sequel', 'joaquin', 'phoenix', "'s", 'joker'], None), ([ 'rt', 'done', 'movies', 'bts', 'bring', 'soul', 'movie', 'we athering', 'joker', 'spiderman', 'far', 'from...'], None)]
```

The data is now ready, we can proceed to building our model.

Section D: Naive Bayes Classifier

Naïve Bayes Classifier is a classification algorithm that relies on Bayes' Theorem, which provides a way of calculating a type of probability called posterior probability, in which the probability of an event A occurring is reliant on probabilistic know background. For example, if person_X only plays tennis when it is not raining outside, then according to Bayesian statistics, the probability of person_X playing tennis when it is not raining can be given as:

$$P(X \text{ plays} \mid \text{no rain}) = P(\text{no rain} \mid X \text{ plays}) * P(x \text{ plays}) / P(\text{no rain})$$

Why Naive Bayes Classifier?

Naive Bayes (NB) are mostly used in natural language processing problems. Naive Bayes predicts the tag of a text by calculating the probability of each tag for a given text and then output the tag with the highest one. It is the most simplest one which is enough in our case.

Other possible classifier

Support Vector Machine:

- effective in high dimensional spaces;
- still effective in cases where number of dimensions is greater than the number of samples
- uses a subset of training points in the decision function (called support vectors), so it is also memory efficient (very important too!)

Conclusion: SVM could be a good option too. Since we need to choose one, we will work with NB.

Step D.1: Build a vocabulary/list of words in our training data set

A vocabulary in NLP is a list of all speech segments available for the model. In our case, this includes all the words present in the Training set. We create a list of distinct words, with its frequency (number of occurrences in the set) as a key, "word_features".

Step D.2: Match tweet content against our vocabulary

I will go through all the words in the Training set, i.e. "word_features", comparing every word against the tweet at hand, associating a number with the word:

- 1 (true): if word in vocabulary is present in tweet
- 0 (false): if word in vocabulary is not resident in tweet

Every word in the vocabulary "word-features", will have a key "contains word X", where X is the word, and associated value "True/False" according to if word in vocabulary is present in tweet or not.

Step D.3: Build our word feature vector

Calling the last two function created, we will build our final feature vector, called "trainingFeatures", with which we can proceed on to training. Our vocabulary of words contains 2700 words.

```
In [40]: trainingFeatures
```

```
Out[40]: [{ 'contains(get)': True, 'contains(swype)': True, 'contains(iphone)': True, 'contains(crack)': True, 'contains(hilarious)': False, 'contains(video)': False, 'contains(guy)': False, 'contains(duet)': False, 'contains(s)': False, 'contains(siri)': False, 'contains(pretty)': False, 'contains(much)': False, 'contains(sums)': False, 'contains(love)': False, 'contains(affair)': False, 'contains(made)': False, 'contains(easy)': False, 'contains(switch)': False, 'contains(see)': False, 'contains(ya)': False, 'contains(16)': False, 'contains(strangest)': False, 'contains(things)': False, 'contains(said)': False, 'contains(far)': False, 'contains(sooo)': False, 'contains(glad)': False, 'contains(gave)': False, 'contains(sense)': False, 'contains(humor)': False, 'contains(via)': False, 'contains(great)': False, 'contains(close)': False, 'contains(personal)': False, 'contains(event)': False, 'contains(tonight)': False, 'contains(regent)': False, 'contains(st)': False, 'contains(store)': False, 'contains(companies)': False, 'contains(experience)': False, 'contains(best)': False
```

Step D.4: Training the classifier

Thanks to the library *nltk* it will take us a short time to train the model as a Naïve Bayes Classifier.

We will perform a small test of our classifier by creating fake sequences of words that are particularly positive and negative to check that the classifier is efficient.

```
test_tweet = ['rt', 'hilarious', 'video', 'guy', 'duet', "'s']
print(NBayesClassifier.classify(extract_features(test_tweet)))
positive

test_tweet = ['terrible', 'bad', 'awful', 'negative', 'duet']
print(NBayesClassifier.classify(extract_features(test_tweet)))
negative
```

Our small and extreme tests show that the classifier is doing some job at least!

Section E: Testing the model

Now, I will run the classifier on the Test set (100 tweets), getting the labels (pos/neg) and calculating the positive or negative percentage of the tweets. That's it!

```
Overall Negative Sentiment  
Negative Sentiment Percentage = 52.0%
```

The results of the sentiment analysis on Tweets about the movie Joker with Joaquin Phoenix that came out in October 2019 is that there is the overall sentiment is negative at 52%. These classification results coincide with the critics one can read in the newspaper and on the internet, as they are balanced between positive and negative. Let's have a deeper look at the results.

Section F: Analyzing the performance of the model

Let's have a look at a couple of Tweets and the assigned label to check if the classification is intuitive or not. This is a first approximate idea of our classifier's performance.

Step F.1: Concatenate pre-processed Tweets and labels

preprocessedTestSet (100, 2)			NBResultLabels (100,1)	
	Tweets pre-processed	label		label
0	[finally, saw, joker, technical, aspects, movi...	None	0	positive
1	[', joker, '', became, first, r-rated, movie,...	None	1	negative
2	[would, like, see, sequel, joaquin, phoenix, '...	None	2	negative
3	[rt, done, movies, bts, bring, soul, movie, we...	None	3	positive
4	[ashishchanchalani, ashishchanchalanimemes, ol...	None	4	negative
5	[think, joker, new, favorite, character/movie,...	None	5	positive
6	[rt, still, thinking, sana, harley, quinn, wan...	None	6	negative
7	[joker, action, movie]	None	7	positive
8	[saw, ford, v, ferrari, last, weekend, great, ...	None	8	positive
9	[', dope, joker, movie]	None	9	positive

```
test_labels
array([[list(['finally', 'saw', 'joker', 'technical', 'as
pects', 'movie', 'stunning', 'acting', 'incredible', 'mig
h...']),
      None, 'positive'],
      [list(['', 'joker', '', 'became', 'first', 'r-
rated', 'movie', 'hit', '1', 'billion', 'box-office', 'ar
e\u200b', 'rumors', 'sequel']),
      None, 'negative'],
      [list(['would', 'like', 'see', 'sequel', 'joaquin'
, 'phoenix', 's', 'joker']),
      None, 'negative'],
      [list(['rt', 'done', 'movies', 'bts', 'bring', 'so
ul', 'movie', 'weathering', 'joker', 'spiderman', 'far',
'from...']),
      None, 'positive'],
      [list(['ashishchanchalani', 'ashishchanchalanimeme
s', 'oldvideos', 'joker', 'think', 'downloaded', 'wrong',
'joker', 'movie']),
      None, 'negative'],
      [list(['think', 'downloaded', 'wrong', 'joker', 'movie'])],
      None, 'negative']])
```


Step F.2: Check 10 pre-processed Tweets and their attributed label

```
# Check first 10 Tweets words and their attributed labels
test_labels[0]
```

```
array([list(['finally', 'saw', 'joker', 'technical', 'aspects', 'movie', 'stunning', 'acting', 'incredible', 'might...']),
      None, 'positive'], dtype=object)
```

```
test_labels[1]
```

```
array([list(['', 'joker', '', 'became', 'first', 'rated', 'movie', 'hit', '1', 'billion', 'box-office', 'are\u200b', 'rumors', 'sequel']),
      None, 'negative'], dtype=object)
```

```
test_labels[2]
```

```
array([list(['would', 'like', 'see', 'sequel', 'joaquin', 'phoenix', 's', 'joker']),
      None, 'negative'], dtype=object)
```

```
test_labels[3]
```

```
array([list(['rt', 'done', 'movies', 'bts', 'bring', 'soul', 'movie', 'weathering', 'joker', 'spiderman', 'far', 'from...']),
      None, 'positive'], dtype=object)
```

```
test_labels[4]
```

```
array([list(['ashishchanchalani', 'ashishchanchalanimemes', 'oldvideos', 'joker', 'think', 'downloaded', 'wrong', 'joker', 'movie']),
      None, 'negative'], dtype=object)
```

```
test_labels[5]
```

```
array([list(['think', 'joker', 'new', 'favorite', 'character/movie', 'love', 'every', 'actor', '', 'played', 'get', 'new', 'parts', 'him...']),
      None, 'positive'], dtype=object)
```

```
test_labels[6]
```

```
array([list(['rt', 'still', 'thinking', 'sana', 'harley', 'quinn', 'wanted', 'chaeyoung', 'joker', 'coz', 'chaeyoung', 'love', 'movie...']),
      None, 'negative'], dtype=object)
```

```
test_labels[7]
```

```
array([list(['joker', 'action', 'movie']), None, 'positive'], dtype=object)
```



```
test_labels[8]
```

```
array([list(['saw', 'ford', 'v', 'ferrari', 'last', 'week  
end', 'great', 'movie', 'would', 'top', 'list', '2019', "  
'm", 'marvel', 'guy', 'avengers...']),  
       None, 'positive'], dtype=object)
```

```
test_labels[9]
```

```
array([list(['', 'dope', 'joker', 'movie']), None, 'posi  
tive'],  
       dtype=object)
```

Step F.3: Model performance

Tweet[0]: At a glance, we see words such as 'stunning', 'incredible' that are positive. The label is therefore 'positive' as expected. CORRECT

Tweet[1]: This tweet seems rather neutral, but is classified as 'negative'. BIAS?

Tweet[2]: This tweet seems rather neutral, but is classified as 'negative'. BIAS?

Tweet[3]: This tweet seems rather neutral, but is classified as 'positive'. BIAS?

Tweet[4]: This tweet seems rather neutral, but is classified as 'negative'. BIAS?

Tweet[5]: This tweet seems rather positive with the word 'love', 'favorite'. The label is therefore 'positive' as expected. CORRECT

Tweet[6]: This tweet seems rather neutral, with a small positive touch with the word 'love'. However, the label is 'negative'. WRONG

Tweet[7]: This tweet seems rather neutral, but is classified as 'negative'. BIAS?

Tweet[8]: At a glance, we see words such as 'great', 'top' that are positive. The label is therefore 'positive' as expected. CORRECT

Tweet[9]: This tweet seems rather neutral, with a small positive note 'dope'. The label is therefore 'positive' as expected. CORRECT

Step F.4: Wrap up observations

1. **50%** | It seems like neutral Tweets are more often classified as negative (4/5) which is a consequence of simplifying the classification options to two (positive and negative). More rarely (1/5), a neutral Tweet is classified as positive.
2. **40%** | Correct classification of the Tweet
3. **10%** | Wrong classification of the Tweet

This (subjective) performance evaluation task is telling us that the classification model is not optimal yet.

Step F.5: Possible improvements

- Trying different combinations of pre-processing steps
- Adding some pre-processing steps
- Checking the impact of each pre-processing step on model performance
- Building our own Training Dataset, to control quality of the classification in the Training Dataset, as it is crucial to a good model performance.

Conclusion

Sentiment Analysis is an interesting way to think about the applicability of Natural Language Processing in making automated conclusions about text. It is being utilized in social media trend analysis and, sometimes, for marketing purposes. Making a Sentiment Analysis program in Python is not a difficult task, thanks to modern-day, ready-for-use libraries.

With this project, we wanted to analyze Tweets about the late released movie "Joker" that attracted all kinds of comments, critics on the web. In summary, we have

- created a Test set directly from Twitter thanks to the API
- used an existing Training set that was classified into positive and negative tweets according to its content by hand
- cleaned both sets by removing any signs (punctuation, hashtags...) that do not bring anything to the sentiment analysis
- created a vocabulary of words based on the Training Set
- matched tweets against vocabulary
- trained the classifier
- tested the model on the test Set
- evaluated the performance of our model

The result of the sentiment analysis is that almost half of the Tweets in the Test set are classified as negative and the other half as positive. The results show the divided critics about the movie, which is a feeling I had when scanning through articles on the web.

However, the quality of our model seems compromised, and several areas of improvement have been identified above.

I would not use the model per se, but most probably after implementing the different changes mentioned above.

Another area of improvement to the model is in regards to measuring the performance of our model. The way the project was conducted didn't allow to use a solid model performance tool such as a confusion matrix. What could be done to measure better the model's performance would be to:

- separate the Training Tweets Dataset between Train and Test sets
- develop and train the model on the Train set
- use Test set to test model
- run metrics (confusion matrix) to evaluate model's performance over Test set and Train set