# Mukesh Chapagain Blog

PHP Magento Nodejs Python Machine Learning Programming & Tutorial

HOME　　ABOUT　　CONTACT　　ADVERTISE　　ARCHIVES　　PRIVACY POLICY

Search the archive...

Home » Machine Learning, Natural Language Processing (NLP), Python, Sentiment Analysis

**3 April 2018**

# Python NLTK: Sentiment Analysis on Movie Reviews [Natural Language Processing (NLP)]

Facebook　Twitter　LinkedIn　Print　Email

This article shows how you can perform sentiment analysis on movie reviews using Python and Natural Language Toolkit (NLTK).

Sentiment Analysis means analyzing the sentiment of a given text or document and categorizing the text/document into a specific class or category (like positive and negative). In other words, we can say that sentiment analysis classifies any particular text or document as positive or negative. Basically, the classification is done for two classes: positive and negative. However, we can add more classes like neutral, highly positive, highly negative, etc.

Sentiment Analysis is also referred as Opinion Mining. It's mostly used in social media and customer reviews data.

In this article, we will learn about labeling data, extracting features, training classifier, and testing the accuracy of the classifier.

## Supervised Classification

Here, we will be doing supervised text classification. In supervised classification, the classifier is trained with **labeled training data**.

In this article, we will use the NLTK's `movie_reviews` corpus as our labeled training data. The `movie_reviews` corpus contains 2K movie reviews with sentiment polarity classification. It's compiled by Pang, Lee.

Here, we have two categories for classification. They are: positive and negative. The movie_reviews corpus already has the reviews categorized as positive and negative.

```
1  from nltk.corpus import movie_reviews
2
3  # Total reviews
4  print (len(movie_reviews.fileids())) # Output: 2000
5
6  # Review categories
7  print (movie_reviews.categories()) # Output: [u'neg', u'pos']
8
9  # Total positive reviews
10 print (len(movie_reviews.fileids('pos'))) # Output: 1000
11
12 # Total negative reviews
13 print (len(movie_reviews.fileids('neg'))) # Output: 1000
14
15 positive_review_file = movie_reviews.fileids('pos')[0]
16 print (positive_review_file) # Output: pos/cv000_29590.txt
```

**Create list of movie review document**

This list contains array containing tuples of all movie review words and their respective category (pos or neg).

```
1  documents = []
2
3  for category in movie_reviews.categories():
4      for fileid in movie_reviews.fileids(category):
5          #documents.append((list(movie_reviews.words(fileid)), category))
6          documents.append((movie_reviews.words(fileid), category))
7
8  print (len(documents)) # Output: 2000
9
10 # x = [str(item) for item in documents[0][0]]
11 # print (x)
```

## About

Mukesh Chapagain is a graduate of Kathmandu University (Dhulikhel, Nepal) from where he holds a Masters degree in Computer Engineering. Mukesh is a passionate web developer who has keen interest in open source technologies, programming & blogging. **more...**

### Subscribe via Email

Email: [_____] Subscribe

### Categories

Select Category

### Archives

Select Month

### Recent Posts

- Magento 2: Create Widget Programmatically & Assign Static Block to it
- Magento 2: Add/Update CMS Static Block via Install/Upgrade Script Programmatically
- Magento 2: Add/Update CMS Page via Install/Upgrade Script Programmatically
- FFMPEG: Convert & Edit Video via Command Line
- ImageMagick: Convert/Edit Multiple Images
- [INFOGRAPHIC] Magento Basic Facts for Newbies
- ImageMagick: Convert & Edit Image via Command Line
- [SOLVED] ERROR 2006 (HY000): MySQL server has gone away
- [SOLVED] MySQL: The server quit without updating PID file
- Magento 2: Change Increment ID Prefix, Suffix, Start value, Step, Pad length of Order, Invoice, Creditmemo & Shipment

### Related Posts

- Python NLTK: Twitter Sentiment Analysis [Natural Language Processing (NLP)]
- Python NLTK: Text Classification [Natural Language Processing (NLP)]
- Natural Language Processing (NLP): Basic Introduction to NLTK [Python]
- Python NLTK: Stop Words [Natural Language Processing (NLP)]
- Python NLTK: Stemming & Lemmatization [Natural Language Processing (NLP)]
- Python NLTK: Working with WordNet [Natural Language Processing (NLP)]
- Python NLTK: Part-of-Speech (POS) Tagging [Natural Language Processing (NLP)]
- Machine Learning & Sentiment Analysis: Text Classification using Python & NLTK
- Python: Twitter Sentiment Analysis on Real Time Tweets using TextBlob
- Python: Twitter Sentiment Analysis using TextBlob

### Most Viewed

- How to Calculate Inverter & Battery Backup Time?
- Very Simple Add, Edit, Delete, View (CRUD) in PHP & MySQL [Beginner Tutorial]
- Magento: How to get attribute name and value?
- CodeIgniter: Simple Add, Edit, Delete, View – MVC CRUD Application
- Magento: How to get controller, module, action and router name?
- Magento: How to select, insert, update, and delete
- LaTeX: Generate dummy text (lorem ipsum) in your document
- CRUD with Login & Register in PHP & MySQL (Add, Delete, View)

Privacy - Terms

```
12
13   # print first tuple
14   print (documents[0])
15   '''
16   Output:
17
18   (['plot', ':', 'two', 'teen', 'couples', 'go', ...], 'neg')
19   '''
20
21   # shuffle the document list
22   from random import shuffle
23   shuffle(documents)
```

## Feature Extraction

To classify the text into any category, we need to define some criteria. On the basis of those criteria, our classifier will learn that a particular kind of text falls in a particular category. This kind of criteria is known as `feature`. We can define one or more feature to train our classifier.

In this example, we will use the `top-N words feature`.

**Fetch all words from the movie reviews corpus**

We first fetch all the words from all the movie reviews and create a list.

```
1   all_words = [word.lower() for word in movie_reviews.words()]
2
3   # print first 10 words
4   print (all_words[:10])
5   '''
6   Output:
7
8   ['plot', ':', 'two', 'teen', 'couples', 'go', 'to', 'a', 'church', 'party']
9   '''
```

**Create Frequency Distribution of all words**

Frequency Distribution will calculate the number of occurence of each word in the entire list of words.

```
1    from nltk import FreqDist
2
3    all_words_frequency = FreqDist(all_words)
4
5    print (all_words_frequency)
6    '''
7    Output:
8
9    <FreqDist with 39768 samples and 1583820 outcomes>
10   '''
11
12   # print 10 most frequently occurring words
13   print (all_words_frequency.most_common(10))
14   '''
15   Output:
16
17   [(',', 77717), ('the', 76529), ('.', 65876), ('a', 38106), ('and', 35576),
     ('of', 34123), ('to', 31937), ("'", 30585), ('is', 25195), ('in', 21822)]
18   '''
```

**Removing Punctuation and Stopwords**

From the above frequency distribution of words, we can see the most frequently occurring words are either punctuation marks or stopwords.

Stop words are those frequently words which do not carry any significant meaning in text analysis. For example, I, me, my, the, a, and, is, are, he, she, we, etc.

Punctuation marks like comma, fullstop. inverted comma, etc. occur highly in any text data.

We will do `data cleaning` by removing stop words and punctuations.

**Remove Stop Words**

```
1    from nltk.corpus import stopwords
2
3    stopwords_english = stopwords.words('english')
4    print (stopwords_english)
5    '''
6    Output:
7
8    ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'your',
     'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'he
```

```
        r', 'hers', 'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'thei
        rs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', 'these',
         'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
         'has', 'had', 'having', 'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'an
        d', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by',
         'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
         'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'ou
        t', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'her
        e', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'f
        ew', 'more', 'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'o
        wn', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just',
         'don', 'should', 'now', 'd', 'll', 'm', 'o', 're', 've', 'y', 'ain', 'are
        n', 'couldn', 'didn', 'doesn', 'hadn', 'hasn', 'haven', 'isn', 'ma', 'might
        n', 'mustn', 'needn', 'shan', 'shouldn', 'wasn', 'weren', 'won', 'wouldn']
 9      '''
10
11      # create a new list of words by removing stopwords from all_words
12      all_words_without_stopwords = [word for word in all_words if word not in sto
        pwords_english]
13
14      # print the first 10 words
15      print (all_words_without_stopwords[:10])
16      '''
17      Output:
18
19      ['plot', ':', 'two', 'teen', 'couples', 'go', 'church', 'party', ',', 'drin
        k']
20      '''
21
22      '''
23      # Above code is written using the List Comprehension feature of Python
24      # It's the same thing as writing the following, the output is the same
25
26      all_words_without_stopwords = []
27      for word in all_words:
28          if word not in stopwords_english:
29              all_words_without_stopwords.append(word)
30
31      print (all_words_without_stopwords[:10])
32      '''
```

You can see that after removing stopwords, the words `to` and `a` has been removed from the first 10 words result.

**Remove Punctuation**

```
 1      import string
 2
 3      print (string.punctuation)
 4      '''
 5      Output:
 6
 7      !"#$%&'()*+,-./:;<=>?@[\]^_`{|}~
 8      '''
 9
10      # create a new list of words by removing punctuation from all_words
11      all_words_without_punctuation = [word for word in all_words if word not in s
        tring.punctuation]
12
13      # print the first 10 words
14      print (all_words_without_punctuation[:10])
15      '''
16      Output:
17
18      [u'plot', u'two', u'teen', u'couples', u'go', u'to', u'a', u'church', u'part
        y', u'drink']
19      '''
```

You can see that on the list that all punctuations like semi-colon `:`, comma `,` are removed.

**Remove both Stopwords & Punctuation**

In the above examples, at first, we only removed stopwords and then in the next code, we only removed punctuation.

Below, we will remove both stopwords and punctuation from the `all_words` list.

```
 1      # Let's name the new list as all_words_clean
 2      # because we clean stopwords and punctuations from the word list
 3
 4      all_words_clean = []
 5      for word in all_words:
 6          if word not in stopwords_english and word not in string.punctuation:
 7              all_words_clean.append(word)
 8
 9      print (all_words_clean[:10])
```

```
10  '''
11  Output:
12
13  ['plot', 'two', 'teen', 'couples', 'go', 'church', 'party', 'drink', 'driv
    e', 'get']
14  '''
```

**Frequency Distribution of cleaned words list**

Below is the frequency distribution of the new list after removing stopwords and punctuation.

```
 1  all_words_frequency = FreqDist(all_words_clean)
 2
 3  print (all_words_frequency)
 4  '''
 5  Output:
 6
 7  <FreqDist with 39586 samples and 710578 outcomes>
 8  '''
 9
10  # print 10 most frequently occurring words
11  print (all_words_frequency.most_common(10))
12  '''
13  Output:
14
15  [('film', 9517), ('one', 5852), ('movie', 5771), ('like', 3690), ('even', 25
    65), ('time', 2411), ('good', 2411), ('story', 2169), ('would', 2109), ('muc
    h', 2049)]
16  '''
```

Previously, before removing stopwords and punctuation, the frequency distribution was:

`FreqDist with 39768 samples and 1583820 outcomes`

Now, the frequency distribution is:

`FreqDist with 39586 samples and 710578 outcomes`

This shows that after removing around 200 stop words and punctuation, the outcomes/words number has reduced to around half of the original size.

The `most common words` or highly occurring words list has also got meaningful words in the list. Before, the first 10 frequently occurring words were only stop-words and punctuations.

**Create Word Feature using 2000 most frequently occurring words**

We take 2000 most frequently occurring words as our feature.

```
 1  print (len(all_words_frequency)) # Output: 39586
 2
 3  # get 2000 frequently occuring words
 4  most_common_words = all_words_frequency.most_common(2000)
 5  print (most_common_words[:10])
 6  '''
 7  Output:
 8
 9  [('film', 9517), ('one', 5852), ('movie', 5771), ('like', 3690), ('even', 25
    65), ('time', 2411), ('good', 2411), ('story', 2169), ('would', 2109), ('muc
    h', 2049)]
10  '''
11
12  print (most_common_words[1990:])
13  '''
14  Output:
15
16  [('genuinely', 64), ('path', 64), ('eve', 64), ('aware', 64), ('bank', 64),
    ('bound', 64), ('eric', 64), ('regular', 64), ('las', 64), ('niro', 64)]
17  '''
18
19  # the most common words list's elements are in the form of tuple
20  # get only the first element of each tuple of the word list
21  word_features = [item[0] for item in most_common_words]
22  print (word_features[:10])
23  '''
24  Output:
25
26  ['film', 'one', 'movie', 'like', 'even', 'time', 'good', 'story', 'would',
    'much']
27  '''
```

**Create Feature Set**

Now, we write a function that will be used to create feature set. The feature set is used to train the classifier.

We define a feature extractor function that checks if the words in a given document are present in the word_features list or not.

```
1  def document_features(document):
2      # "set" function will remove repeated/duplicate tokens in the given list
3      document_words = set(document)
4      features = {}
5      for word in word_features:
6          features['contains(%s)' % word] = (word in document_words)
7      return features
8
9  # get the first negative movie review file
10 movie_review_file = movie_reviews.fileids('neg')[0]
11 print (movie_review_file)
12 '''
13 Output:
14
15 neg/cv000_29416.txt
16 '''
17
18 #print (document_features(movie_reviews.words(movie_review_file)))
19 '''
20 Output:
21
22 {'contains(waste)': False, 'contains(lot)': False, 'contains(rent)': False,
    'contains(black)': False, 'contains(rated)': False, 'contains(potential)': F
    alse,
23 ..............................................................................
   ................................ 'contains(smile)': False, 'con
24 tains(cross)': False, 'contains(barry)': False}
   '''
```

In the beginning of this article, we have created the `documents` list which contains data of all the movie reviews. Its elements are tuples with word list as first item and review category as the second item of the tuple.

```
1  # print first tuple of the documents list
2  print (documents[0])
3  '''
4  Output:
5
6  (['plot', ':', 'two', 'teen', 'couples', 'go', ...], 'neg')
7  '''
```

We now loop through the `documents` list and create a feature set list using the `document_features` function defined above.

– Each item of the feature_set list is a tuple.

– The first item of the tuple is the dictionary returned from `document_features` function

– The second item of the tuple is the category (pos or neg) of the movie review

```
1  feature_set = [(document_features(doc), category) for (doc, category) in doc
   uments]
2  print (feature_set[0])
3  '''
4  Output:
5
6  ({'contains(waste)': False, 'contains(lot)': False, 'contains(rent)': False,
   'contains(black)': False, 'contains(rated)': False, 'contains(potential)': F
   alse,
   ..............................................................................
7  ................................................. 'contains(goo
   d)': False, 'contains(live)': False, 'contains(synopsis)': False, 'contains
   (appropriate)': False, 'contains(towards)': False, 'contains(smile)': False,
   'contains(cross)': False, 'contains(barry)': False}, 'neg')
8  '''
9
10 '''
11 # In the above code, we have used list-comprehension feature of python
12 # The same code can be written as below:
13 feature_set = []
14 for (doc, category) in documents:
15     feature_set.append((document_features(doc), category))
16 print (feature_set[0])
17 '''
```

## Training Classifier

From the feature set we created above, we now create a separate training set and a separate testing/validation set. The train set is used to train the classifier and the test set is used to test the classifier to check how accurately it classifies the given text.

### Creating Train and Test Dataset

In this example, we use the first 400 elements of the feature set array as a test set and the rest of the data as a train set. Generally, 80/20 percent is a fair split between training and testing set, i.e. 80 percent training set and 20 percent testing set.

```python
1  print (len(feature_set)) # Output: 2000
2
3  test_set = feature_set[:400]
4  train_set = feature_set[400:]
5
6  print (len(train_set)) # Output: 1600
7  print (len(test_set)) # Output: 400
```

### Training a Classifier

Now, we train a classifier using the training dataset. There are different kind of classifiers namely Naive Bayes Classifier, Maximum Entropy Classifier, Decision Tree Classifier, Support Vector Machine Classifier, etc.

In this example, we use the Naive Bayes Classifier. It's a simple, fast, and easy classifier which performs well for small datasets. It's a simple probabilistic classifier based on applying Bayes' theorem. Bayes' theorem describes the probability of an event, based on prior knowledge of conditions that might be related to the event.

```python
1  from nltk import NaiveBayesClassifier
2
3  classifier = NaiveBayesClassifier.train(train_set)
```

### Testing the trained Classifier

Let's see the accuracy percentage of the trained classifier. The accuracy value changes each time you run the program because of the names array being shuffled above.

```python
1  from nltk import classify
2
3  accuracy = classify.accuracy(classifier, test_set)
4  print (accuracy) # Output: 0.77
```

Let's see the output of the classifier by providing some custom reviews.

```python
1  from nltk.tokenize import word_tokenize
2
3  custom_review = "I hated the film. It was a disaster. Poor direction, bad ac
   ting."
4  custom_review_tokens = word_tokenize(custom_review)
5  custom_review_set = document_features(custom_review_tokens)
6  print (classifier.classify(custom_review_set)) # Output: neg
7  # Negative review correctly classified as negative
8
9  # probability result
10 prob_result = classifier.prob_classify(custom_review_set)
11 print (prob_result) # Output: <ProbDist with 2 samples>
12 print (prob_result.max()) # Output: neg
13 print (prob_result.prob("neg")) # Output: 0.999989264571
14 print (prob_result.prob("pos")) # Output: 1.07354285262e-05
15
16 custom_review = "It was a wonderful and amazing movie. I loved it. Best dire
   ction, good acting."
17 custom_review_tokens = word_tokenize(custom_review)
18 custom_review_set = document_features(custom_review_tokens)
19
20 print (classifier.classify(custom_review_set)) # Output: neg
21 # Positive review is classified as negative
22 # We need to improve our feature set for more accurate prediction
23
24 # probability result
25 prob_result = classifier.prob_classify(custom_review_set)
26 print (prob_result) # Output: <ProbDist with 2 samples>
27 print (prob_result.max()) # Output: neg
28 print (prob_result.prob("neg")) # Output: 0.999791868552
29 print (prob_result.prob("pos")) # Output: 0.000208131447797
```

Let's see the most informative features among the entire features in the feature set.

```
1  # show 5 most informative features
2  print (classifier.show_most_informative_features(10))
3  '''
4  Output:
5
6  Most Informative Features
7     contains(outstanding) = True              pos : neg    =     14.7 : 1.0
8          contains(mulan) = True              pos : neg    =      7.8 : 1.0
9         contains(poorly) = True              neg : pos    =      7.7 : 1.0
10   contains(wonderfully) = True              pos : neg    =      7.5 : 1.0
11         contains(seagal) = True              neg : pos    =      6.5 : 1.0
12          contains(awful) = True              neg : pos    =      6.1 : 1.0
13         contains(wasted) = True              neg : pos    =      6.1 : 1.0
14          contains(waste) = True              neg : pos    =      5.6 : 1.0
15          contains(damon) = True              pos : neg    =      5.3 : 1.0
16          contains(flynt) = True              pos : neg    =      5.1 : 1.0
17  '''
```

The result shows that the word `outstanding` is used in positive reviews 14.7 times more often than it is used in negative reviews the word `poorly` is used in negative reviews 7.7 times more often than it is used in positive reviews. Similarly, for other letters. These ratios are also called **likelihood ratios**.

Therefore, a review has a high chance to be classified as positive if it contains words like `outstanding` and `wonderfully`. Similarly, a review has a high chance of being classified as negative if it contains words like `poorly`, `awful`, `waste`, etc.

**Note:** You can modify the `document_features` function to generate the feature set which can improve the accuracy of the trained classifier. Feature extractors are built through a process of trail-and-error & guided by intuitions.

## Bag of Words Feature

In the above example, we used **top-N words** feature. We used 2000 most frequently occurring words as our top-N words feature. The classifier identified negative review as negative. However, the classifier was not able to classify positive review correctly. It classified a positive review as negative.

**Top-N words feature**

– The top-N words feature is also a bag-of-words feature.
– But in the top-N feature, we only used the top 2000 words in the feature set.
– We combined the positive and negative reviews into a single list, randomized the list, and then separated the train and test set.
– This approach can result in the un-even distribution of positive and negative reviews across the train and test set.

**Bag-of-words feature shown below**

In the bag-of-words feature as shown below:

– We will use all the useful words of each review while creating the feature set.
– We take a fixed number of positive and negative reviews for train and test set.
– This result in equal distribution of positive and negative reviews across train and test set.

In the approach shown below, we will modify the feature extractor function.

– We form a list of unique words of each review.
– The category (pos or neg) is assigned to each bag of words.
– Then the category of any given text is calculated by matching the different bag-of-words & their respective category.

```
1  from nltk.corpus import movie_reviews
2
3  pos_reviews = []
4  for fileid in movie_reviews.fileids('pos'):
5      words = movie_reviews.words(fileid)
6      pos_reviews.append(words)
7
8  neg_reviews = []
9  for fileid in movie_reviews.fileids('neg'):
10     words = movie_reviews.words(fileid)
11     neg_reviews.append(words)
12
13 # print first positive review item from the pos_reviews list
14 print (pos_reviews[0])
15 '''
16 Output:
```

```
17
18   ['films', 'adapted', 'from', 'comic', 'books', ...]
19   '''
20
21   # print first negative review item from the neg_reviews list
22   print (neg_reviews[0])
23   '''
24   Output:
25
26   ['plot', ':', 'two', 'teen', 'couples', 'go', ...]
27   '''
28
29   # print first 20 items of the first item of positive review
30   print (pos_reviews[0][:20])
31   '''
32   Output:
33
34   ['films', 'adapted', 'from', 'comic', 'books', 'have', 'had', 'plenty', 'o
     f', 'success', ',', 'whether', 'they', "'", 're', 'about', 'superheroes',
      '(', 'batman', ',']
35   '''
36
37   # print first 20 items of the first item of negative review
38   print (neg_reviews[0][:20])
39   '''
40   Output:
41
42   ['plot', ':', 'two', 'teen', 'couples', 'go', 'to', 'a', 'church', 'party',
      ',', 'drink', 'and', 'then', 'drive', '.', 'they', 'get', 'into', 'an']
43   '''
```

### Feature Extraction

We use the bag-of-words feature. Here, we clean the word list (i.e. remove stop words and punctuation). Then, we create a dictionary of cleaned words.

```
1    from nltk.corpus import stopwords
2    import string
3
4    stopwords_english = stopwords.words('english')
5
6    # feature extractor function
7    def bag_of_words(words):
8        words_clean = []
9
10       for word in words:
11           word = word.lower()
12           if word not in stopwords_english and word not in string.punctuation:
13               words_clean.append(word)
14
15       words_dictionary = dict([word, True] for word in words_clean)
16
17       return words_dictionary
18
19   # using dict will remove duplicate words from the words list
20   # note the output: stopword 'the' is also removed
21   print (bag_of_words(['the', 'the', 'good', 'bad', 'the', 'good']))
22   '''
23   Output:
24
25   {'bad': True, 'good': True}
26   '''
```

### Create Feature Set

We use the bag-of-words feature and tag each review with its respective category as positive or negative.

```
1    # positive reviews feature set
2    pos_reviews_set = []
3    for words in pos_reviews:
4        pos_reviews_set.append((bag_of_words(words), 'pos'))
5
6    # negative reviews feature set
7    neg_reviews_set = []
8    for words in neg_reviews:
9        neg_reviews_set.append((bag_of_words(words), 'neg'))
10
11   # print first positive review item from the pos_reviews list
12   print (pos_reviews_set[0])
13   '''
14   Output:
15
16   ({'childs': True, 'steve': True, 'surgical': True, 'go': True, 'certainly': Tr
     True, 'song': True, 'simpsons': True, 'novel': True,
      ......................................................................
17
```

```
18  ..................................................... 'menace': True, 'sta
19  riginal': True}, 'pos')
    '''
20
21
22  # print first negative review item from the neg_reviews list
23  print (neg_reviews_set[0])
    '''
24
25  Output:

    ({'concept': True, 'skip': True, 'insight': True, 'playing': True, 'executed'
26  ue, 'still': True, 'find': True, 'seemed': True,

27  .......................................................................
    ...................................................... 'entertaining': True, 'years
    y': True, 'came': True}, 'neg')
    '''
```

**Create Train and Test Set**

There are 1000 positive reviews set and 1000 negative reviews set. We take 20% (i.e. 200) of positive
reviews and 20% (i.e. 200) of negative reviews as a test set. The remaining negative and positive
reviews will be taken as a training set.

> **Note:**
> – There is difference between pos_reviews & pos_reviews_set array which are defined above.
> – pos_reviews array contains words list only
> – pos_reviews_set array contains words feature list
> – pos_reviews_set & neg_reviews_set arrays are used to create train and test set as shown below

```
1   print (len(pos_reviews_set), len(neg_reviews_set)) # Output: (1000, 1000)
2
3   # radomize pos_reviews_set and neg_reviews_set
4   # doing so will output different accuracy result everytime we run the progra
    m
5   from random import shuffle
6   shuffle(pos_reviews_set)
7   shuffle(neg_reviews_set)
8
9   test_set = pos_reviews_set[:200] + neg_reviews_set[:200]
10  train_set = pos_reviews_set[200:] + neg_reviews_set[200:]
11
12  print(len(test_set),  len(train_set)) # Output: (400, 1600)
```

**Training Classifier and Calculating Accuracy**

We train Naive Bayes Classifier using the training set and calculate the classification accuracy of the
trained classifier using the test set.

```
1   from nltk import classify
2   from nltk import NaiveBayesClassifier
3
4   classifier = NaiveBayesClassifier.train(train_set)
5
6   accuracy = classify.accuracy(classifier, test_set)
7   print(accuracy) # Output: 0.7325
8
9   print (classifier.show_most_informative_features(10))
10  '''
11  Output:
12
13  Most Informative Features
14          breathtaking = True              pos : neg    =    20.3 : 1.0
15             dazzling = True               pos : neg    =    12.3 : 1.0
16            ludicrous = True               neg : pos    =    12.2 : 1.0
17          outstanding = True               pos : neg    =    10.6 : 1.0
18             insipid = True                neg : pos    =    10.3 : 1.0
19            stretched = True               neg : pos    =    10.3 : 1.0
20            stupidity = True               neg : pos    =    10.2 : 1.0
21              annual = True                pos : neg    =     9.7 : 1.0
22            headache = True                neg : pos    =     9.7 : 1.0
23              avoids = True                pos : neg    =     9.7 : 1.0
24  '''
```

**Testing Classifier with Custom Review**

We provide custom review text and check the classification output of the trained classifier. The
classifier correctly predicts both negative and positive reviews provided.

```
1   from nltk.tokenize import word_tokenize
2
3
```

Privacy - Terms

```
 4   custom_review = "I hated the film. It was a disaster. Poor direction, bad ac
     ting."
 5   custom_review_tokens = word_tokenize(custom_review)
 6   custom_review_set = bag_of_words(custom_review_tokens)
 7   print (classifier.classify(custom_review_set)) # Output: neg
 8   # Negative review correctly classified as negative
 9
10   # probability result
11   prob_result = classifier.prob_classify(custom_review_set)
12   print (prob_result) # Output: <ProbDist with 2 samples>
13   print (prob_result.max()) # Output: neg
14   print (prob_result.prob("neg")) # Output: 0.776128854994
15   print (prob_result.prob("pos")) # Output: 0.223871145006
16
17   custom_review = "It was a wonderful and amazing movie. I loved it. Best dire
     ction, good acting."
18   custom_review_tokens = word_tokenize(custom_review)
19   custom_review_set = bag_of_words(custom_review_tokens)
20
21   print (classifier.classify(custom_review_set)) # Output: pos
22   # Positive review correctly classified as positive
23
24   # probability result
25   prob_result = classifier.prob_classify(custom_review_set)
26   print (prob_result) # Output: <ProbDist with 2 samples>
27   print (prob_result.max()) # Output: pos
28   print (prob_result.prob("neg")) # Output: 0.0972171562901
     print (prob_result.prob("pos")) # Output: 0.90278284371
```

## Bi-gram Features

N-grams are common terms in text processing and analysis. N-grams are related with words of a text.
There are different n-grams like unigram, bigram, trigram, etc.

> Unigram = Item having a single word, i.e. the n-gram of size 1. For example, good.
>
> Bigram = Item having two words, i.e. the n-gram of size 2. For example, very good.
>
> Trigram = Item having three words, i.e. the n-gram of size 3. For example, not so good.

In the above bag-of-words model, we only used the unigram feature. In the example below, we will use both unigram and bigram feature, i.e. we will deal with both single words and double words.

**Feature Extraction**

In this case, both unigrams and bigrams are used as features.

We define two functions:

– **bag_of_words**: that extracts only unigram features from the movie review words
– **bag_of_ngrams**: that extracts only bigram features from the movie review words

We then define another function:

– **bag_of_all_words**: that combines both unigram and bigram features

```
 1   from nltk import ngrams
 2   from nltk.corpus import stopwords
 3   import string
 4
 5   stopwords_english = stopwords.words('english')
 6
 7   # clean words, i.e. remove stopwords and punctuation
 8   def clean_words(words, stopwords_english):
 9       words_clean = []
10       for word in words:
11           word = word.lower()
12           if word not in stopwords_english and word not in string.punctuation
     :
13               words_clean.append(word)
14       return words_clean
15
16   # feature extractor function for unigram
17   def bag_of_words(words):
18       words_dictionary = dict([word, True] for word in words)
19       return words_dictionary
20
21   # feature extractor function for ngrams (bigram)
22   def bag_of_ngrams(words, n=2):
23       words_ng = []
24       for item in iter(ngrams(words, n)):
25           words_ng.append(item)
26       words_dictionary = dict([word, True] for word in words_ng)
27       return words_dictionary
28
29   '''
```

```python
30  # Alternative Bi-gram feature extractor
31  # using BigramCollocationFinder module
32
33  # Collocations are multiple words which commonly co-occur.
34  # http://www.nltk.org/howto/collocations.html
35  # https://streamhacker.com/2010/05/24/text-classification-sentiment-analysi
    s-stopwords-collocations/
36
37  import itertools
38  from nltk.collocations import BigramCollocationFinder
39  from nltk.metrics import BigramAssocMeasures
40
41  # feature extractor function for ngrams (bigram)
42  # get 200 most frequently occurring bigrams from every review
43  def bag_of_ngrams(words, score_fn=BigramAssocMeasures.chi_sq, n=200):
44      bigram_finder = BigramCollocationFinder.from_words(words)
45      bigrams = bigram_finder.nbest(score_fn, n)
46      return dict([(ngram, True) for ngram in itertools.chain(words, bigram
    s)])
47  '''
48
49  from nltk.tokenize import word_tokenize
50  text = "It was a very good movie."
51  words = word_tokenize(text.lower())
52
53  print (words)
54  '''
55  Output:
56
57  ['it', 'was', 'a', 'very', 'good', 'movie', '.']
58  '''
59
60  print (bag_of_ngrams(words))
61  '''
62  Output:
63
64  {('very', 'good'): True, ('movie', '.'): True, ('it', 'was'): True, ('goo
    d', 'movie'): True, ('was', 'a'): True, ('a', 'very'): True}
65  '''
66
67  # working with cleaning words
68  # i.e. removing stopwords and punctuation
69  words_clean = clean_words(words, stopwords_english)
70  print (words_clean)
71  '''
72  Output:
73
74  ['good', 'movie']
75  '''
76
77  # cleaning words is find for unigrams
78  # but this can omit important words for bigrams
79  # for example, stopwords like very, over, under, so, etc. are important for
    bigrams
80  # we create a new stopwords list specifically for bigrams by omitting such
    important words
81  important_words = ['above', 'below', 'off', 'over', 'under', 'more', 'most'
    , 'such', 'no', 'nor', 'not', 'only', 'so', 'than', 'too', 'very', 'just',
    'but']
82
83  stopwords_english_for_bigrams = set(stopwords_english) - set(important_word
    s)
84
85  words_clean_for_bigrams = clean_words(words, stopwords_english_for_bigrams)
86  print (words_clean_for_bigrams)
87  '''
88  Output:
89
90  ['very', 'good', 'movie']
91  '''
92
93  # We will use general stopwords for unigrams
94  # And special stopwords list for bigrams
95  unigram_features = bag_of_words(words_clean)
96  print (unigram_features)
97  '''
98  Output:
99
100 {'movie': True, 'good': True}
101 '''
102
103 bigram_features = bag_of_ngrams(words_clean_for_bigrams)
104 print (bigram_features)
105 '''
106 Output:
107
108 {('very', 'good'): True, ('good', 'movie'): True}
109 '''
110
111 # combine both unigram and bigram features
112 all_features = unigram_features.copy()
113 all_features.update(bigram_features)
114 print (all_features)
```

```
115  '''
116  Output:
117
118  {'movie': True, ('very', 'good'): True, 'good': True, ('good', 'movie'): Tr
     ue}
119  '''
120
121  # let's define a new function that extracts all features
122  # i.e. that extracts both unigram and bigrams features
123  def bag_of_all_words(words, n=2):
124      words_clean = clean_words(words, stopwords_english)
125      words_clean_for_bigrams = clean_words(words, stopwords_english_for_bigr
     ams)
126
127      unigram_features = bag_of_words(words_clean)
128      bigram_features = bag_of_ngrams(words_clean_for_bigrams)
129
130      all_features = unigram_features.copy()
131      all_features.update(bigram_features)
132
133      return all_features
134
135  print (bag_of_all_words(words))
136  '''
137  Output:
138
139  {'movie': True, ('very', 'good'): True, 'good': True, ('good', 'movie'): Tr
     ue}
140  '''
```

Working with NLTK's movie reviews corpus

```
1   from nltk.corpus import movie_reviews
2
3   pos_reviews = []
4   for fileid in movie_reviews.fileids('pos'):
5       words = movie_reviews.words(fileid)
6       pos_reviews.append(words)
7
8   neg_reviews = []
9   for fileid in movie_reviews.fileids('neg'):
10      words = movie_reviews.words(fileid)
11      neg_reviews.append(words)
```

**Create Feature Set**

```
1   # positive reviews feature set
2   pos_reviews_set = []
3   for words in pos_reviews:
4       pos_reviews_set.append((bag_of_all_words(words), 'pos'))
5
6   # negative reviews feature set
7   neg_reviews_set = []
8   for words in neg_reviews:
9       neg_reviews_set.append((bag_of_all_words(words), 'neg'))
```

**Create Train and Test Set**

There are 1000 positive reviews set and 1000 negative reviews set. We take 20% (i.e. 200) of positive reviews and 20% (i.e. 200) of negative reviews as the test set. The remaining negative and positive reviews will be taken as the training set.

```
1   print (len(pos_reviews_set), len(neg_reviews_set)) # Output: (1000, 1000)
2
3   # radomize pos_reviews_set and neg_reviews_set
4   # doing so will output different accuracy result everytime we run the progra
    m
5   from random import shuffle
6   shuffle(pos_reviews_set)
7   shuffle(neg_reviews_set)
8
9   test_set = pos_reviews_set[:200] + neg_reviews_set[:200]
10  train_set = pos_reviews_set[200:] + neg_reviews_set[200:]
11
12  print(len(test_set),  len(train_set)) # Output: (400, 1600)
```

**Training Classifier and Calculating Accuracy**

We train Naive Bayes Classifier using the training set and calculate the classification accuracy of the trained classifier using the test set.

```
1   from nltk import classify
2   from nltk import NaiveBayesClassifier
```

```
3
4  classifier = NaiveBayesClassifier.train(train_set)
5
6  accuracy = classify.accuracy(classifier, test_set)
7  print(accuracy) # Output: 0.8025
8
9  print (classifier.show_most_informative_features(10))
10 '''
11 Output:
12
13 Most Informative Features
14              insulting = True          neg : pos    =    17.0 : 1.0
15            outstanding = True          pos : neg    =    14.7 : 1.0
16         ('nice', 'see') = True         pos : neg    =    11.7 : 1.0
17        ('one', 'worst') = True         neg : pos    =    11.4 : 1.0
18      ('would', 'think') = True         neg : pos    =    11.0 : 1.0
19       ('quite', 'well') = True         pos : neg    =    11.0 : 1.0
20         ('makes', 'no') = True         neg : pos    =    10.3 : 1.0
21       ('but', 'script') = True         neg : pos    =    10.3 : 1.0
22    ('quite', 'frankly') = True         neg : pos    =    10.3 : 1.0
23              animators = True          pos : neg    =    10.3 : 1.0
24 '''
```

**Note:**

– The accuracy of the classifier has significantly increased when trained with combined feature set (unigram + bigram).

– Accuracy was 73% while using only Unigram features.

– Accuracy has increased to 80% while using combined (unigram + bigram) features.

**Testing Classifier with Custom Review**

We provide custom review text and check the classification output of the trained classifier. The classifier correctly predicts both negative and positive reviews provided.

```
1  from nltk.tokenize import word_tokenize
2
3  custom_review = "I hated the film. It was a disaster. Poor direction, bad acting."
4  custom_review_tokens = word_tokenize(custom_review)
5  custom_review_set = bag_of_all_words(custom_review_tokens)
6  print (classifier.classify(custom_review_set)) # Output: neg
7  # Negative review correctly classified as negative
8
9  # probability result
10 prob_result = classifier.prob_classify(custom_review_set)
11 print (prob_result) # Output: <ProbDist with 2 samples>
12 print (prob_result.max()) # Output: neg
13 print (prob_result.prob("neg")) # Output: 0.770612685688
14 print (prob_result.prob("pos")) # Output: 0.229387314312
15
16
17 custom_review = "It was a wonderful and amazing movie. I loved it. Best direction, good acting."
18 custom_review_tokens = word_tokenize(custom_review)
19 custom_review_set = bag_of_all_words(custom_review_tokens)
20
21 print (classifier.classify(custom_review_set)) # Output: pos
22 # Positive review correctly classified as positive
23
24 # probability result
25 prob_result = classifier.prob_classify(custom_review_set)
26 print (prob_result) # Output: <ProbDist with 2 samples>
27 print (prob_result.max()) # Output: pos
28 print (prob_result.prob("neg")) # Output: 0.00677736186354
29 print (prob_result.prob("pos")) # Output: 0.993222638136
```

References:

1. Learning to Classify Text
2. From Text Classification to Sentiment Analysis

Hope this helps. Thanks.

**Share this:**

Facebook    Twitter    LinkedIn    Print    Email

**Related Posts:**

1. **Python NLTK: Twitter Sentiment Analysis [Natural Language Processing (NLP)]**

2. **Python NLTK: Text Classification [Natural Language Processing (NLP)]**
3. **Natural Language Processing (NLP): Basic Introduction to NLTK [Python]**
4. **Python NLTK: Stop Words [Natural Language Processing (NLP)]**
5. **Python NLTK: Stemming & Lemmatization [Natural Language Processing (NLP)]**
6. **Python NLTK: Working with WordNet [Natural Language Processing (NLP)]**
7. **Python NLTK: Part-of-Speech (POS) Tagging [Natural Language Processing (NLP)]**
8. **Machine Learning & Sentiment Analysis: Text Classification using Python & NLTK**
9. **Python: Twitter Sentiment Analysis on Real Time Tweets using TextBlob**
10. **Python: Twitter Sentiment Analysis using TextBlob**

---

**Get New Post by Email**

[_____]   [Subscribe]

**Find me on**

[f] [t] [g+] [in] [rss]

---

**ALSO ON MUKESH CHAPAGAIN'S BLOG**

**Node.js, MongoDB & Express: Simple …**

3 years ago • 1 comment

This article shows how you can create a simple CRUD (Create, Read, Update, …

**Magento 2: Upload File & Image**

2 years ago • 1 comment

This article shows how you can upload files and images in Magento 2. You can …

**SSH/SCP: Copy Files Folders from Local …**

2 years ago • 3 comments

This article shows how you can copy files and folders from your local machine …

---

♡ **Recommend**  1          🐦 **Tweet**          f **Share**          Sort by Best ▾

0 Comments     **Mukesh Chapagain's Blog**     🔒 **Disqus' Privacy Policy**     **Login** ▾

1

Start the discussion…

LOG IN WITH                OR SIGN UP WITH DISQUS ?

[ Name ]

Be the first to comment.

---

**recent posts**

Magento 2: Create Widget Programmatically & Assign Static Block to it

Magento 2: Add/Update CMS Static Block via Install/Upgrade Script Programmatically

Magento 2: Add/Update CMS Page via Install/Upgrade Script Programmatically

FFMPEG: Convert & Edit Video via Command Line

ImageMagick: Convert/Edit Multiple Images

**most viewed**

How to Calculate Inverter & Battery Backup Time?

Very Simple Add, Edit, Delete, View (CRUD) in PHP & MySQL [Beginner Tutorial]

Magento: How to get attribute name and value?

CodeIgniter: Simple Add, Edit, Delete, View – MVC CRUD Application

Magento: How to get controller, module, action and router name?

**tag cloud**

admin array attribute category checkout command line configurable product crud currency customer Database error extension football Google grid HTML image Javascript jQuery latex Linux login machine learning Magento magento2 module MySQL natural language processing NLP nltk nodejs order payment PHP product python query shopping cart template Twitter Ubuntu Wordpress XML

Privacy - Terms

Privacy - Terms