

03/07/2020

Feature Extraction with frequencies

→ Feature Extraction

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

↓ ↓ ↓

features Bias Sum of Pos.
of tweet'm' frequencies

↓

Sum of Neg.
frequencies

e.g.

→ I am sad, I am not learning NLP

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

→ Now, if we look at the table - word frequencies
in classes

then the data on which we have to look

is

vocabulary	Pos. freq.(1)	Neg. freq.
I	3	3
am	3	3

we don't have
to look at
these because

in our example

these words

are not

there

Learning	1	1
NLP	1	1
sad	0	2
not	0	1

→ So,

$$\sum_w \text{freats}(w, 1) = 8$$

$$\sum_w \text{freats}(w, 0) = 11$$

$$\rightarrow x_m = [1, 8, 11]$$

Preprocessing: Stop words and punctuation

Eg: → @Yousri and @AndrewNG are tuning

the GREAT AI model at <https://deeplearning.ai>!!!

<u>Stop words</u>	<u>Punctuation</u>
and	,
is	.
are	:
at	!
has	"
for	'
a	

→ Now, tweet will be

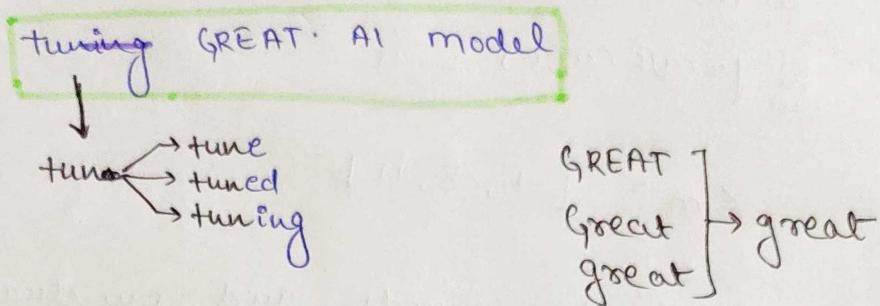
@Yousri and @AndrewNG tuning GREAT AI model
<https://deeplearning.ai>

Preprocessing: Handles & URLs

→ @Yousri and @AndrewNG tuning GREAT AI model
<https://deeplearning.ai>

⇒ tuning GREAT AI model

Preprocessing: stemming and lowercasing



Preprocessed tweet:

[tun, great, ai, model]

EXERCISE

→ for this exercise, we will use a twitter dataset that comes with NLTK.

→ About dataset

- The sample dataset from NLTK is separated into positive & negative tweets.
- It contains 5000 positive tweets and 5000 negative tweets exactly.

Import basic libraries

```
import nltk  
from nltk.corpus import  
import matplotlib.pyplot as plt  
import random
```

→ Python library for NLP

twitter_samples → sample twitter dataset from NLTK.

→ Library for visualization

→ pseudo-random number generators.

- * To download twitter dataset at your local machine you can type:
 - `nltk.download("twitter_samples")`
- # we can load the text fields of positive & negative tweets by using the module's `strings()` method like this:
 - `all-positive-tweets = twitter_samples.strings('positive_tweets.json')`
 - `all-negative-tweets = twitter_samples.strings('negative_tweets.json')`
 - `print('Number of positive tweets:', len(all-positive-tweets))`
↳ 5000
 - `print('Number of negative tweets:', len(all-negative-tweets))`
↳ 5000
↳ <class 'list'>
 - `print('The type of all-positive tweets is:', type(all-negative-tweets[0]))`
↳ <class 'str'>
- # Declare a figure with a custom size


```
fig = plt.figure(figsize=(5,5))
```
- # Labels for the two classes


```
labels = 'Positive', 'Negative'
```
- # Sizes for each slice


```
sizes = [len(all-positive-tweets), len(all-negative-tweets)]
```
- # Declare pie chart, where the slices will be ordered and plotted counter-clockwise:


```
plt.pie(sizes, labels=labels, autopct='%1.1f%%',  
        shadow=True, startangle=90)
```

```
# Equal aspect ratio ensures that pie is drawn as a circle.
```

```
plt.axis('equal')
```

```
# Display the chart
```

```
plt.show()
```



→ # Print random positive & negative tweet in green & red.

```
print('\033[92m' + all_positive_tweets[random.randint(0,500)])
```

```
print('\033[91m' + all_negative_tweets[random.randint(0,500)])
```

Output:

i walked out of the fruit shop with lots of pumpkins & a guy smiled at me because it looked pretty funny, but doesn't matter, it worked :)
really didn't want to see that :-c

```
# Preprocess raw text for sentiment analysis
```

→ The preprocessing steps are comprised of the following tasks:

- Tokenizing the string
- Lowercasing
- Removing stop words & punctuation
- Stemming

→ # Our selected sample. Complex enough to exemplify each step.

```
tweet = all_positive_tweets[2277]
print(tweet)
```

→ # Let's import few more library for this purpose.

```
# download the stopwords from NLTK
nltk.download('stopwords')
```

→ import re → library for regular expression operations
import string → for string operations
from nltk.corpus import stopwords → module for stop words that comes with NLTK.
from nltk.stem import PorterStemmer → module for stemming
from nltk.tokenize import TweetTokenizer → module for tokenizing strings

remove old style retweet text "RT"

```
tweet2 = re.sub(r'^RT[\S]*', '', tweet)
```

remove hyperlinks

```
tweet2 = re.sub(r'https?:\/\/.*[\r\n]', '', tweet2)
```

remove hashtags (only remove the hash # sign from the word)

```
tweet2 = re.sub(r'\#', '', tweet2)
```

```
print(tweet2)
```

Tokenize the string

→ # Instantiate tokenizer class

```
tokenizer = TweetTokenizer(preserve_case=False,  
                           strip_handles=True, reduce_len=True)
```

Tokenize tweets

```
tweet_tokens = tokenizer.tokenize(tweet2)
```

```
print('Tokenized string:\n' + tweet_tokens)
```

Remove stop words and punctuations

→ # Import the english stop words list from NLTK

```
stopwords_english = stopwords.words('english')
```

```
print('Stop words\n')
```

```
print(stopwords_english)
```

```
print('\nPunctuation\n')
```

```
print(string.punctuation)
```

→ tweets_clean = []

```
for word in tweet_tokens: # Go through every word  
    if (word not in stopwords_english and # remove  
        word not in string.punctuation): # remove  
            tweets_clean.append(word) # punctuation
```

```
print('Removed stop words & punctuation:')
```

```
print(tweets_clean)
```

Stemming

→ # Instantiate stemming class

stemmer = PorterStemmer()

Create an empty list to store the stems

tweets_stem = []

for word in tweets_clean:

stem_word = stemmer.stem(word) # Stemming word

tweets_stem.append(stem_word) # Append to the list

print('stemmed words:')

print(tweets_stem)

process_tweet()

→ As shown above, preprocessing consists of multiple steps before you arrive at the final list of words.

→ There is a function "process_tweet(tweet)" available in "util.py" which performs all the same preprocessing steps which is discussed above.

→ from util import process_tweet # Import process_tweet function

choose the same tweet

tweet = all_positive_tweets[2277]

print('1033[92m' + tweet)

print('1033[94m')

call the imported function

tweets_stem = process_tweet(tweet); # Preprocess a given tweet

print('Preprocessed tweet:')

print(tweets_stem) # Print the result.