

# Text Processing

One of the biggest breakthroughs required for achieving any level of artificial intelligence is to have machines which can process text data. Thankfully, the amount of text data being generated in this universe has exploded exponentially in the last few years.

It has become imperative for an organization to have a structure in place to mine actionable insights from the text being generated. From social media analytics to risk management and cybercrime protection, dealing with text data has never been more important.

Here we will discuss different feature extraction methods, starting with some basic techniques and we will also learn about pre-processing of the text data in order to extract better features from clean data.

## 1. Basic Feature Extraction

In the tutorial, we will use a twitter sentiment dataset its available on Moodle.

```
import pandas as pd

train = pd.read_csv('train_E6oV3lV.csv')
```

*Note that here we are only working with textual data, but we can also use the below methods when numerical features are also present along with the text.*

### 1.1 Number of Words

One of the most basic features we can extract is the number of words in each tweet. The basic intuition behind this is that generally, the negative sentiments contain a lesser amount of words than the positive ones.

To do this, we simply use the **split** function in python:

```
train['word_count'] = train['tweet'].apply(lambda x: len(str(x).split(" ")))
train[['tweet', 'word_count']].head()
```

	tweet	word_count
0	@user when a father is dysfunctional and is s...	21
1	@user @user thanks for #lyft credit i can't us...	22
2	bihday your majesty	5
3	#model i love u take with u all the time in ...	17
4	factsguide: society now #motivation	8

## 1.2 Number of characters

This feature is also based on the previous feature intuition. Here, we calculate the number of characters in each tweet. This is done by calculating the length of the tweet.

```
train['char_count'] = train['tweet'].str.len()
## this also includes spaces
train[['tweet', 'char_count']].head()
```

	tweet	char_count
0	@user when a father is dysfunctional and is s...	102
1	@user @user thanks for #lyft credit i can't us...	122
2	bihday your majesty	21
3	#model i love u take with u all the time in ...	86
4	factsguide: society now #motivation	39

Note that the calculation will also include the number of spaces, which you can remove, if required.

## 1.3 Average Word Length

We will also extract another feature which will calculate the average word length of each tweet. This can also potentially help us in improving our model.

Here, we simply take the sum of the length of all the words and divide it by the total length of the tweet:

```
def avg_word(sentence):
    words = sentence.split()
    return (sum(len(word) for word in words)/len(words))

train['avg_word'] = train['tweet'].apply(lambda x: avg_word(x))
train[['tweet', 'avg_word']].head()
```

	tweet	avg_word
0	@user when a father is dysfunctional and is s...	4.555556
1	@user @user thanks for #lyft credit i can't us...	5.315789
2	bihday your majesty	5.666667
3	#model i love u take with u all the time in ...	4.928571
4	factsguide: society now #motivation	8.000000

## 1.4 Number of stopwords

Generally, while solving an NLP problem, the first thing we do is to remove the stopwords. But sometimes calculating the number of stopwords can also give us some extra information which we might have been losing before.

Here, we have imported stopwords from *NLTK*, which is a basic NLP library in python.

```
import nltk

nltk.download('stopwords')

from nltk.corpus import stopwords
stop = stopwords.words('english')

train['stopwords'] = train['tweet'].apply(lambda x: len([x for x in
x.split() if x in stop]))
train[['tweet', 'stopwords']].head()
```

	tweet	stopwords
0	@user when a father is dysfunctional and is s...	10
1	@user @user thanks for #lyft credit i can't us...	4
2	bihday your majesty	1
3	#model i love u take with u all the time in ...	5
4	factsguide: society now #motivation	1

## 1.5 Number of special characters

One more interesting feature which we can extract from a tweet is calculating the number of hashtags or mentions present in it. This also helps in extracting extra information from our text data.

Here, we make use of the *'starts with'* function because hashtags (or mentions) always appear at the beginning of a word.

```
train['hashtags'] = train['tweet'].apply(lambda x: len([x for x in x.split()
if x.startswith('#')]))
train[['tweet', 'hashtags']].head()
```

	tweet	hashtags
0	@user when a father is dysfunctional and is s...	1
1	@user @user thanks for #lyft credit i can't us...	3
2	bihday your majesty	0
3	#model i love u take with u all the time in ...	1
4	factsguide: society now #motivation	1

## 1.6 Number of numerics

Just like we calculated the number of words, we can also calculate the number of numerics which are present in the tweets. It does not have a lot of use in our example, but this is still a useful feature that should be run while doing similar exercises. For example,

```
train['numerics'] = train['tweet'].apply(lambda x: len([x for x in
x.split() if x.isdigit()])))
train[['tweet', 'numerics']].head()
```

	tweet	numerics
0	@user when a father is dysfunctional and is s...	0
1	@user @user thanks for #lyft credit i can't us...	0
2	bihday your majesty	0
3	#model i love u take with u all the time in ...	0
4	factsguide: society now #motivation	0

## 1.7 Number of Uppercase words

Anger or rage is quite often expressed by writing in UPPERCASE words which makes this a necessary operation to identify those words.

```
train['upper'] = train['tweet'].apply(lambda x: len([x for x in x.split()
if x.isupper()])))
train[['tweet', 'upper']].head()
```

	tweet	upper
0	@user when a father is dysfunctional and is s...	0
1	@user @user thanks for #lyft credit i can't us...	0
2	bihday your majesty	0
3	#model i love u take with u all the time in ...	0
4	factsguide: society now #motivation	0

## 2. Basic Pre-processing

So far, we have learned how to extract basic features from text data. Before diving into text and feature extraction, our first step should be cleaning the data in order to obtain better features. We will achieve this by doing some of the basic pre-processing steps on our training data.

So, let's get into it.

## 2.1 Lower case

The first pre-processing step which we will do is transform our tweets into lower case. This avoids having multiple copies of the same words. For example, while calculating the word count, ‘Analytics’ and ‘analytics’ will be taken as different words.

```
train['tweet'] = train['tweet'].apply(lambda x: " ".join(x.lower() for x in
x.split()))
train['tweet'].head()
0    @user when a father is dysfunctional and is so...
1    @user @user thanks for #lyft credit i can't us...
2                                     bihday your majesty
3    #model i love u take with u all the time in ur...
4                                     factsguide: society now #motivation
Name: tweet, dtype: object
```

## 2.2 Removing Punctuation

The next step is to remove punctuation, as it doesn't add any extra information while treating text data. Therefore removing all instances of it will help us reduce the size of the training data.

```
train['tweet'] = train['tweet'].str.replace('[^\w\s]','')
train['tweet'].head()
0    user when a father is dysfunctional and is so ...
1    user user thanks for lyft credit i cant use ca...
2                                     bihday your majesty
3    model i love u take with u all the time in urð...
4                                     factsguide society now motivation
Name: tweet, dtype: object
```

As you can see in the above output, all the punctuation, including ‘#’ and ‘@’, has been removed from the training data.

## 2.3 Removal of Stop Words

As we discussed earlier, stop words (or commonly occurring words) should be removed from the text data. For this purpose, we can either create a list of stopwords ourselves or we can use predefined libraries.

```
from nltk.corpus import stopwords
stop = stopwords.words('english')
train['tweet'] = train['tweet'].apply(lambda x: " ".join(x for x in
x.split() if x not in stop))
train['tweet'].head()
```

```

0    user father dysfunctional selfish drags kids d...
1    user user thanks lyft credit cant use cause do...
2                                bihday majesty
3    model love u take u time urð ðððð ððð
4    factsguide society motivation
Name: tweet, dtype: object

```

## 2.4 Common word removal

Previously, we just removed commonly occurring words in a general sense. We can also remove commonly occurring words from our text data. First, let's check the 10 most frequently occurring words in our text data then take call to remove or retain.

```

freq = pd.Series(' '.join(train['tweet']).split()).value_counts()[:10]
freq
> user      17473
love       2647
ð          2511
day        2199
â          1797
happy      1663
amp        1582
im         1139
u          1136
time       1110
dtype: int64

```

Now, let's remove these words as their presence will not of any use in classification of our text data.

```

freq = list(freq.index)
train['tweet'] = train['tweet'].apply(lambda x: " ".join(x for x in
x.split() if x not in freq))
train['tweet'].head()

```

---

```

0    father dysfunctional selfish drags kids dysfun...
1    thanks lyft credit cant use cause dont offer w...
2                                bihday majesty
3    model take urð ðððð ððð
4    factsguide society motivation
Name: tweet, dtype: object

```

## 2.5 Rare words removal

Similarly, just as we removed the most common words, this time let's remove rarely occurring words from the text. Because they're so rare, the association between them and other words is dominated by noise. You can replace rare words with a more general form and then this will have higher counts

```

freq = pd.Series(' '.join(train['tweet']).split()).value_counts()[-10:]
freq
> tvperfect      1
oau              1

```

```

850am          1
semangatpagi  1
kindestbravest 1
moodyah        1
downhill       1
loreal         1
ohwhatcoulditbe 1
maannnn        1
dtype: int64

```

```

freq = list(freq.index)
train['tweet'] = train['tweet'].apply(lambda x: " ".join(x for x in
x.split() if x not in freq))
train['tweet'].head()

```

```

0    father dysfunctional selfish drags kids dysfun...
1    thanks lyft credit cant use cause dont offer w...
2                                     bihday majesty
3                                     model take urð ðððð ððð
4                                     factsguide society motivation
Name: tweet, dtype: object

```

All these pre-processing steps are essential and help us in reducing our vocabulary clutter so that the features produced in the end are more effective.

## 2.6 Spelling correction

We've all seen tweets with a plethora of spelling mistakes. Our timelines are often filled with hastily sent tweets that are barely legible at times.

In that regard, spelling correction is a useful pre-processing step because this also will help us in reducing multiple copies of words. For example, "Analytics" and "analytcs" will be treated as different words even if they are used in the same sense.

To achieve this we will use the *textblob* library.

```
!pip install -U textblob
```

```

from textblob import TextBlob
train['tweet'][:5].apply(lambda x: str(TextBlob(x).correct()))
0    father dysfunctional selfish drags kiss dysfun...
1    thanks left credit can use cause dont offer wh...
2                                     midday majesty
3                                     model take or ðððð ððð
4                                     factsguide society motivation
Name: tweet, dtype: object

```

Note that it will actually take a lot of time to make these corrections. Therefore, just for the purposes of learning, I have shown this technique by applying it on only the first 5 rows. Moreover, we cannot always expect it to be accurate so some care should be taken before applying it.



We should also keep in mind that words are often used in their abbreviated form. *For instance, 'your' is used as 'ur'.* We should treat this before the spelling correction step, otherwise these words might be transformed into any other word like the one shown below:

```
0 father dysfunctional selfish drags kiss dysfun...
1 thanks left credit can use cause dont offer wh...
2                                     midday majesty
3                                     model take or ôôôô ôôô
4                                     factsguide society motivation
Name: tweet, dtype: object
```

## 2.7 Tokenization

Tokenization refers to dividing the text into a sequence of words or sentences. In our example, we have used the *textblob* library to first transform our tweets into a blob and then converted them into a series of words.

```
nltk.download('punkt')

TextBlob(train['tweet'][1]).words
> WordList(['thanks', 'lyft', 'credit', 'cant', 'use', 'cause', 'dont',
'offer', 'wheelchair', 'vans', 'pdx', 'disappointed', 'getthanked'])
```

## 2.8 Stemming

Stemming refers to the removal of suffices, like “ing”, “ly”, “s”, etc. by a simple rule-based approach. For this purpose, we will use *PorterStemmer* from the NLTK library.

```
from nltk.stem import PorterStemmer
st = PorterStemmer()
train['tweet'][:5].apply(lambda x: " ".join([st.stem(word) for word in
x.split()])))
0 father dysfunct selfish drag kid dysfunct run
1 thank lyft credit cant use caus dont offer whe...
2                                     bihday majesti
3                                     model take urô ôôôô ôôô
4                                     factsguid societi motiv
Name: tweet, dtype: object
```

In the above output, *dysfunctional* has been transformed into *dysfunct*, among other changes.

## 2.9 Lemmatization

Lemmatization is a more effective option than stemming because it converts the word into its root word, rather than just stripping the suffices. It makes use of the vocabulary and does a morphological analysis to obtain the root word. Therefore, **we usually prefer using lemmatization over stemming.**

```
from textblob import Word
train['tweet'] = train['tweet'].apply(lambda x: "
".join([Word(word).lemmatize() for word in x.split()])))
train['tweet'].head()
0 father dysfunctional selfish drag kid dysfunct...
```



```

1      thanks lyft credit cant use cause dont offer w...
2                                          bihday majesty
3                                          model take urð öððð öðð
4                                          factsguide society motivation
Name: tweet, dtype: object

```

### 3. Advance Text Processing

Up to this point, we have done all the basic pre-processing steps in order to clean our data. Now, we can finally move on to extracting features using NLP techniques.

#### 3.1 N-grams

N-grams are the combination of multiple words used together. Ngrams with N=1 are called unigrams. Similarly, bigrams (N=2), trigrams (N=3) and so on can also be used.

Unigrams do not usually contain as much information as compared to bigrams and trigrams. The basic principle behind n-grams is that they capture the language structure, like what letter or word is likely to follow the given one. The longer the n-gram (the higher the  $n$ ), the more context you have to work with. Optimum length really depends on the application – if your n-grams are too short, you may fail to capture important differences. On the other hand, if they are too long, you may fail to capture the “general knowledge” and only stick to particular cases.

So, let’s quickly extract bigrams from our tweets using the *ngrams* function of the *textblob* library.

```

TextBlob(train['tweet'][0]).ngrams(2)
> [WordList(['user', 'when']),
   WordList(['when', 'a']),
   WordList(['a', 'father']),
   WordList(['father', 'is']),
   WordList(['is', 'dysfunctional']),
   WordList(['dysfunctional', 'and']),
   WordList(['and', 'is']),
   WordList(['is', 'so']),
   WordList(['so', 'selfish']),
   WordList(['selfish', 'he']),
   WordList(['he', 'drags']),
   WordList(['drags', 'his']),
   WordList(['his', 'kids']),
   WordList(['kids', 'into']),
   WordList(['into', 'his']),
   WordList(['his', 'dysfunction']),
   WordList(['dysfunction', 'run'])]

```

### 3.2 Term frequency

Term frequency is simply the ratio of the count of a word present in a sentence, to the length of the sentence.

Therefore, we can generalize term frequency as:

**TF = (Number of times term T appears in the particular row) / (number of terms in that row)**

Below, I have tried to show you the term frequency table of a tweet.

```
tfl = (train['tweet'][1:2]).apply(lambda x: pd.value_counts(x.split("
"))).sum(axis = 0).reset_index()
tfl.columns = ['words', 'tf']
tfl
```

	words	tf
0	lyft	1
1	offer	1
2	wheelchair	1
3	van	1
4	pdx	1
5	cause	1
6	cant	1
7	use	1
8	thanks	1
9	credit	1
10	getthanked	1
11	disappointed	1
12	dont	1

### 3.3 Inverse Document Frequency

The intuition behind inverse document frequency (IDF) is that a word is not of much use to us if it's appearing in all the documents.

Therefore, the IDF of each word is the log of the ratio of the total number of rows to the number of rows in which that word is present.

**IDF =  $\log(N/n)$** , where, N is the total number of rows and n is the number of rows in which the word was present.

So, let's calculate IDF for the same tweets for which we calculated the term frequency.

```
import numpy as np
for i,word in enumerate(tf1['words']):
    tf1.loc[i, 'idf'] =
np.log(train.shape[0]/(len(train[train['tweet'].str.contains(word)])))

tf1
```

	words	tf	idf
0	lyft	1	8.762865
1	offer	1	6.522155
2	wheelchair	1	9.273691
3	van	1	5.236505
4	pdx	1	8.762865
5	cause	1	5.690172
6	cant	1	3.538194
7	use	1	3.552287
8	thanks	1	4.597751
9	credit	1	7.327781
10	getthanked	1	9.679156
11	disapointed	1	10.372303
12	dont	1	3.745585

The more the value of IDF, the more unique is the word.

### 3.4 Term Frequency – Inverse Document Frequency (TF-IDF)

TF-IDF is the multiplication of the TF and IDF which we calculated above.

```
tf1['tfidf'] = tf1['tf'] * tf1['idf']
tf1
```

	words	tf	idf	tfidf
0	lyft	1	8.762865	8.762865
1	offer	1	6.522155	6.522155
2	wheelchair	1	9.273691	9.273691
3	van	1	5.236505	5.236505
4	pdx	1	8.762865	8.762865
5	cause	1	5.690172	5.690172
6	cant	1	3.538194	3.538194
7	use	1	3.552287	3.552287
8	thanks	1	4.597751	4.597751
9	credit	1	7.327781	7.327781
10	getthanked	1	9.679156	9.679156
11	disappointed	1	10.372303	10.372303
12	dont	1	3.745585	3.745585

We can see that the TF-IDF has penalized words like ‘don’t’, ‘can’t’, and ‘use’ because they are commonly occurring words. However, it has given a high weight to “disappointed” since that will be very useful in determining the sentiment of the tweet.

We don’t have to calculate TF and IDF every time beforehand and then multiply it to obtain TF-IDF. Instead, *sklearn* has a separate function to directly obtain it:

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer(max_features=1000, lowercase=True, analyzer='word',
    stop_words= 'english', ngram_range=(1,1))
train_vect = tfidf.fit_transform(train['tweet'])

train_vect
<31962x1000 sparse matrix of type '<class 'numpy.float64'>'
    with 114033 stored elements in Compressed Sparse Row format>
```

We can also perform basic pre-processing steps like lower-casing and removal of stopwords, if we haven’t done them earlier.

### 3.5 Bag of Words

Bag of Words (BoW) refers to the representation of text which describes the presence of words within the text data. The intuition behind this is that two similar text fields will contain similar kind of words, and will therefore have a similar bag of words. Further, that from the text alone we can learn something about the meaning of the document.

For implementation, *sklearn* provides a separate function for it as shown below:

```

from sklearn.feature_extraction.text import CountVectorizer
bow = CountVectorizer(max_features=1000, lowercase=True,
ngram_range=(1,1), analyzer = "word")
train_bow = bow.fit_transform(train['tweet'])
train_bow
> <31962x1000 sparse matrix of type '<class 'numpy.int64'>'
   with 128380 stored elements in Compressed Sparse Row format>

```

### 3.6 Sentiment Analysis

If you recall, our problem was to detect the sentiment of the tweet. So, before applying any ML/DL models (which can have a separate feature detecting the sentiment using the textblob library), let's check the sentiment of the first few tweets.

```

train['tweet'][:5].apply(lambda x: TextBlob(x).sentiment)
0      (-0.3, 0.5354166666666667)
1      (0.2, 0.2)
2      (0.0, 0.0)
3      (0.0, 0.0)
4      (0.0, 0.0)
Name: tweet, dtype: object

```

Above, you can see that it returns a tuple representing polarity and subjectivity of each tweet. Here, we only extract polarity as it indicates the sentiment as value nearer to 1 means a positive sentiment and values nearer to -1 means a negative sentiment. This can also work as a feature for building a machine learning model.

```

train['sentiment'] = train['tweet'].apply(lambda x:
TextBlob(x).sentiment[0] )
train[['tweet','sentiment']].head()

```

	tweet	sentiment
0	father dysfunctional selfish drag kid dysfunct...	-0.3
1	thanks lyft credit cant use cause dont offer w...	0.2
2	bihday majesty	0.0
3	model take urð ðððð ððð	0.0
4	factsguide society motivation	0.0

## End Notes

I hope that now you have a basic understanding of how to deal with text data in predictive modelling. These methods will help in extracting more information which in return will help you in building better models.