



Corpus Length									
Unique Words									
Lexical Richness									

## Task 2: Term Frequency

Similar to lexical richness, if you want to determine the count or frequency of a particular word appearing in a corpus, you can use:

```
text1.count('monster') / len(text1) * 100
```

Define a function **TF()** which can take the corpus and the token and return this value for you. In the case of log scales, you can use:

```
math.log(text1.count('monster')+1,10)
```

For this, use a function **LOGTF()** which can do the same.

For Inverse Document Frequency, you can use:

```
math.log(9 / text1.count('monster'), 10)
```

Prepare a function **IDF()** for the above.

Obtain the 10 most common words in a given corpus using the Frequency Distribution:

```
fdist1 = FreqDist(text1)
fdist1.most_common(3)
```

You would naturally want to use remove some stop words from the corpus (See Task 3). To view the most common words as a plot, use:

```
fdist1.plot(50)
```

Now, fill the table below:

	Text1		
Tokens	TF()	LOGTF()	IDF()
monster			
evil			
devil			
the			
Common word 1			
Common word 2			
Common word 3			

## Task 3: Tokenization & POS

To try it out, try a simple tokenization task using the NLTK Punkt tokenizer.

```
nltk.download('punkt')

text = "NLTK is a powerful library for natural language processing."

words = nltk.word_tokenize(text)

sentences = nltk.sent_tokenize(text)

print(words)
print(sentences)
```

To to the text, assign parts of speech tags to the sentences:

```
nltk.download('averaged_perceptron_tagger')

tags = nltk.pos_tag(words)

print(tags)
```

Remove some of the stop words like is, a, etc. in the sentence:

```
nltk.download('stopwords')

from nltk.corpus import stopwords

filtered_words = [word for word in words if word.lower() not in
stopwords.words('english')]

print(filtered_words)
```