

## Data is imported

Make Sure that cik\_list.xlsx and Python.py are in same folder

```
In [1]: import pandas as pd

df = pd.read_excel('cik_list.xlsx')
```

Complete the url of SECFNAME

```
import urllib, requests, urllib.request
```

## Extract URL text and store

```
list=[]
list_data=[]
for i in df['SECNAME_url']:
```

[illegible]

```
import nltk.data
```

```
nltk.download('punkt')
sentence_list=[]
for i in list_data:
    sentence_list.append(nltk.tokenize.sent_tokenize(i))
```

```
import re
tokenizedsent_list=[]
```

```
word_list=[]
for y in i:
    y=y.strip()
    processed_text = y.lower()
    processed_text = re.sub('[^a-zA-Z]', ' ', processed_text )
    processed_text = re.sub(r'[+]', ' ', processed_text)
    processed_text = re.sub('[<[^\>]+>]', ' ', processed_text)

    all_words = nltk.sent_tokenize(processed_text)

    word_list+=(nltk.word_tokenize(sent) for sent in all_words)
tokenizedsent_list.append(word_list)
```

Make sure all the files are in same folder as Python.py

```
master = pd.read_excel ('LoughranMcDonald_MasterDictionary_201
shortgeneric=open('StopWords_Generic.txt','r').readlines()
```

```
currencies=open('StopWords_Currencies.txt','r').readlines()
auditor=open('StopWords_Auditor.txt','r').readlines()
dataandnumbers=open('StopWords_DatesandNumbers.txt','r').readlines()
geographic=open('StopWords_Geographic.txt','r').readlines()
names=open('StopWords_Names.txt','r').readlines()
```

```
stopwords = []
```

```
for i in shortgeneric, longgeneric, currencies, auditor, dataandnumbers, geographic, names:
    for y in i:
        y = re.partition('|')(0)
        y = y.strip()
        y = y.lower()
        stopwords.append(y)
```

```
positive=[]
for i in range(len(master['Word'])):
    if master['Positive'][i]!=0:
```

```
negative=[]
for i in range(len(master['Word'])):
    if master['Negative'][i]!=0:
        negative.append(master['Word'][i])

complex_words=[]
for i in range(len(master['Word'])):
    if master['Syllables'][i]>2:
        complex_words.append(master['Word'][i].lower())
```

```
constraining=pd.read_excel ('constraining_dictionary.xlsx')['Word'].to_list()
uncertainly=pd.read_excel ('uncertainty_dictionary.xlsx')['Word'].to_list()
```

exclusive a combined dictionary sentiment

```
for i in negative,positive,uncertainly,constraining:
    for element in i:
        sentiment_dict[element.lower()] = s
    s+=2
```

with their count

```
def counter(slist):
```

```
def wordDict():
    for y in list:
        if y not in dict.keys():
            dict[y]=1
        else:
            dict[y]+=1
    return dict
```

```
dictsentence=[]
s=0
for i in range(1000000):
```

```
dummy=[  
    for y in range(len(tokenizedsent_list[i])):  
        dummy+=(tokenizedsent_list[i][y])  
    dictsentence.append(counter(dummy))
```

```
processed_sent=[]
s=0
```

```
dummy=i.copy()
for y in i.keys():
```

```

1         del dummy[y]
2         processed_sent.append(dummy)

```

## Each word in the list

positive,negative,constraint and uncertain words for a document are increased accordingly

```

def f1():
    score_sentiment=[]
    s=0
    for i in processed_sent:
        positive_count=0
        negative_count=0
        complex_score=0
        constraining_count=0
        uncertainty_count=0
        for y in i.keys():
            if y in sentiment_dict.keys():
                if (sentiment_dict[y]==-1):
                    negative_count+=i[y]
                elif (sentiment_dict[y]==1):
                    positive_count+=i[y]
                elif (sentiment_dict[y]==3):
                    uncertainty_count+=i[y]
                elif (sentiment_dict[y]==5):
                    constraining_count+=i[y]
            if y in complex_words:
                complex_score+=i[y]
        s+=1
    score_sentiment.append([positive_count,negative_count,constraining_count,uncertainty_count,complex_score])
    print(s,[positive_count,negative_count,constraining_count,uncertainty_count,complex_score])

```

The total count of words and sentences for each document are calculated

```
total_data=[]
for i in range(len(processed_sent)):
    count=0
    for y in processed_sent[i].keys():
        count+=processed_sent[i][y]
    total_data.append([count,len(tokenizedsent_list[i])])
```

```
positive_score=[]
negative_score=[]
polarity_score=[]
average_sentence_length=[]
```

```
percentage_of_complex_words=[]
fog_index=[]
complex_word_count=[]
word_count=[]
uncertainty_score=[]
constraining_score=[]
positive_word_proportion=[]
negative_word_proportion=[]
```

```
constraining_words_whole_report=[]
count_cwvr=0
for cik in range(len(procesed_sent)):
```

```
cik_n=score_sentiment[cik][1]
cik_c=score_sentiment[cik][2]
cik_u=score_sentiment[cik][3]
cik_cwc=score_sentiment[cik][4]
cik_tw=total_data[cik][0]
cik_ts=total_data[cik][1]
positive_score.append(cik_p)
negative_score.append(cik_n)
constraining_score.append(cik_c)
count_cwvr+=cik_c
uncertainty_score.append(cik_u)
complex_word_count.append(cik_cwc)
word_count.append(cik_tw)
average_sentence_length.append(cik_tw/cik_ts)
percentage_of_complex_words.append(cik_cwc/cik_tw)
fog_index.append(0.4*(cik_tw/cik_ts+cik_cwc/cik_tw))
positive_word_proportion.append(cik_p/cik_tw)
negative_word_proportion.append(cik_n/cik_tw)
constraining_word_proportion.append(cik_c/cik_tw)
uncertainty_word_proportion.append(cik_u/cik_tw)
polarity_score.append((cik_p-cik_n)/(cik_p+cik_n+0.000001))
if (cik==(len(procesed_sent)-1)):
    constraining_words_whole_report.append(count_cwvr)
else:
    constraining_words_whole_report.append('')
```

## The list are exported as Output\_Score.xlsx

```
In [ ]: data = {'CIK':df['CIK'].to_list(),
               'CONAME':df['CONAME'].to_list(),
               'FYRMO':df['FYRMO'].to_list(),
               'FDATE':df['FDATE'].to_list(),
               'FORM':df['FORM'].to_list(),
               'SECFNAME':df['SECFNAME'].to_list(),
               'positive_score':positive_score,
               'negative_score':negative_score,
               'polarity_score':polarity_score,
               'average_sentence_length':average_sentence_length,
```

```
'fog_i
'compl
```

```
'uncertainty_score':uncertainty_score,
'constraining_score':constraining_score,
'positive_word_proportion':positive_word_proportion,
'negative_word_proportion':negative_word_proportion,
'uncertainty_word_proportion':uncertainty_word_proportion,
'constraining_word_proportion':constraining_word_proportion,
'constraining_words_whole_report':constraining_words_whole_report]
```

```
'polarity_score'
'word_count', 'un
```

```
file.to_excel('Output_Score.xlsx', index = False )
```

```
In [ ]: polarity_score
```