MINERAÇÃO DE TEXTOS (Parte 2) - WORD PAIRS

Livro: The Amateur Emigrant **de** Robert Louis Stevenson

Usando o algoritmo de WordPairs, do pacote NETWORKX, para descobrir as associações entre as palavras do texto.

Instalando o pacote 'NETWORKX' e as bibliotecas necessárias

```
In [1]:
         pip install networkx
        Requirement already satisfied: networkx in c:\users\mmateus\anaconda3\lib\site-packa
        ges (2.6.3)
        Note: you may need to restart the kernel to use updated packages.
In [2]:
         import pandas as pd
         import matplotlib.pyplot as plt
         import re
         import nltk
         import networkx as nx
         from nltk import ngrams
         from nltk.tokenize import word_tokenize
         from nltk.corpus import stopwords
         from nltk.corpus import opinion_lexicon
         nltk.download('stopwords', quiet=True)
         nltk.download('opinion_lexicon', quiet=True)
         nltk.download('punkt', quiet=True)
        True
Out[2]:
```

Importação do texto:

```
In [3]:
    df = pd.read_csv('data/TheAmateurEmigrant.txt',sep='\t')\
        .dropna()
    df.head(10)
```

```
Out[3]:

THE AMATEUR EMIGRANT

THE SECOND CABIN

I first encountered my fellow-passengers on th...

Thence we descended the Clyde in no familiar s...

on each other as on possible enemies. A few S...

laready grown acquainted on the North Sea, wer...

their long pipes; but among English speakers d...
```

- 7 reigned supreme. The sun was soon overclouded...
- 8 grew sharp as we continued to descend the wide...
- 9 falling temperature the gloom among the passen...

Preparação do texto

```
def clean_text(text):
    text = text.lower()
    text = text.replace("'", '')
    text = re.sub(r'[^\w]', ' ', text)
    text = re.sub(r'\s+', ' ', text)
    text = text.strip()
    return text

df['text'] = df['text'].map(clean_text)
    df['text'] = df['text'].map(word_tokenize)
    df.head()
```

```
Out[4]:

0 [the, amateur, emigrant]

1 [the, second, cabin]

2 [i, first, encountered, my, fellow, passengers...

3 [thence, we, descended, the, clyde, in, no, fa...

4 [on, each, other, as, on, possible, enemies, a...
```

Criando tokens de pares de palavras (WORD PAIRS)

```
In [5]:
    df['wordpairs'] = df['text'].map(lambda x: list(ngrams(x, 2)))
    df = df.explode('wordpairs')

    df.head(10)
```

Out[5]:		text	wordpairs
	0	[the, amateur, emigrant]	(the, amateur)
	0	[the, amateur, emigrant]	(amateur, emigrant)
	1	[the, second, cabin]	(the, second)
	1	[the, second, cabin]	(second, cabin)
	2	[i, first, encountered, my, fellow, passengers	(i, first)
	2	[i, first, encountered, my, fellow, passengers	(first, encountered)
	2	[i, first, encountered, my, fellow, passengers	(encountered, my)
	2	[i, first, encountered, my, fellow, passengers	(my, fellow)
	2	[i, first, encountered, my, fellow, passengers	(fellow, passengers)
	2	[i, first, encountered, my, fellow, passengers	(passengers, on)

Classificando os tokens por frequência

```
In [6]:
          df['wordpairs'].value_counts().head(10)
         (of, the)
                        181
Out[6]:
         (in, the)
                        150
         (to, the)
                         84
         (he, was)
                         81
                         80
         (and, the)
                         72
         (it, was)
         (on, the)
                         63
         (he, had)
                         63
         (of, a)
                         62
         (with, a)
                         53
         Name: wordpairs, dtype: int64
```

Separando as Word Pairs para remover as Stop Words

```
In [7]:
    df = pd.DataFrame(df.wordpairs.values.tolist(), columns=['word1', 'word2']).dropna()
    df.head(10)
```

```
Out[7]:
                   word1
                                 word2
          0
                       the
                                amateur
          1
                  amateur
                                emigrant
          2
                       the
                                 second
          3
                   second
                                   cabin
           4
                         i
                                    first
          5
                      first
                            encountered
              encountered
                                  fellow
          7
                       my
          8
                    fellow
                             passengers
               passengers
                                     on
```

```
In [8]: df.shape
Out[8]: (24030, 2)
```

Removing Stop Words

```
In [9]:
    en_stopwords = set(stopwords.words('english'))
    df = df[~(df.word1.isin(en_stopwords) | df.word2.isin(en_stopwords))]
    df.head()
```

```
Out[9]: word1 word2
```

```
1 amateur emigrant
3 second cabin
5 first encountered
8 fellow passengers
21 familiar spirit
```

```
In [10]: df.shape

Out[10]: (4005, 2)
```

Com a remoção das Stop Words, o texto foi reduzido em 84% de palavras, mantendo apenas as que são relevantes.

Classificando as Word Pairs por frequência

word1	word2	n
second	cabin	18
new	york	12
fellow	passengers	11
sick	man	6
next	morning	5
dare	say	5
could	see	5
fellow	passenger	4
mr	jones	4
one	day	4
hurricane	deck	4
steerage	passenger	4
brass	plate	4
working	man	4
one	night	4
emigrant	ship	3
would	say	3
	second new fellow sick next dare could fellow mr one hurricane steerage brass working one emigrant	second cabin new york fellow passengers sick man next morning dare say could see fellow passenger mr jones one day hurricane deck steerage passenger brass plate working man one night emigrant ship

word1

word2

	word1	word2	n
3589	west	street	3
1837	long	ago	3
335	broken	meat	3

Visualizando os dados através de uma tabela.

• Restrição das word pairs para apenas aquelas que aparecem mais de 3 vezes.

```
In [12]: df[df.n > 3].head(10)
```

Out[12]: word1 word2 n **2872** second cabin 18 2207 york 12 new 1035 fellow passengers 11 2970 sick man 6 2213 next morning 657 dare 5 say 611 could see 1034 fellow passenger

mr

one

jones

day

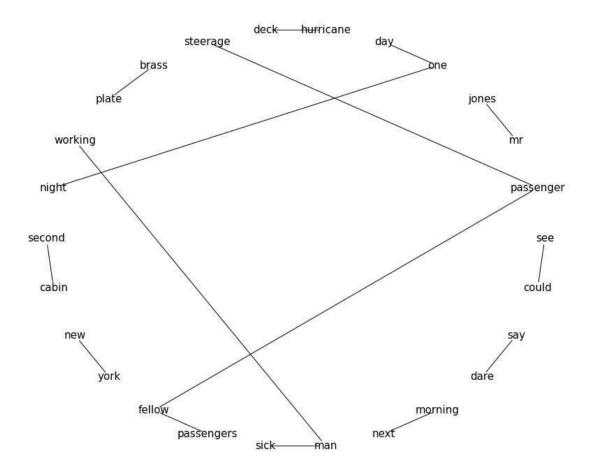
4

2106

2321

Visualizando o gráfico de conexões entre os pares de palavras

```
In [13]:
    G = nx.from_pandas_edgelist(df[df.n > 3], 'word1', 'word2')
    plt.figure(figsize=(12, 10))
    nx.draw_shell(G, with_labels=True, node_color='white', font_size=15)
```



Através da associação de palavras acima, dá para se ter uma ideia de que o livro de Stevenson trata de uma história que se passa num navio, por usar palavras como "steerage passenger" (passageiro da terceira classe), "second cabin" (cabine da segunda classe) e "working man" (trabalhador).