# 2ª Lista de Exercícios - Mineração de Texto

Maria Fernanda Souza Andrade, Saulo Lucas Ferreira

01/10/2018

# 1 Usando o NLTK

## 1.1 Gerar um word cloud contendo os 20 lemas de substantivos mais frequentes.

```
arq = open('NLP.txt','r')
text = arq.read()
arq.close()

lemmatizer = WordNetLemmatizer()
# tokenizing
tokens = nltk.word_tokenize(text)
# lemmatizing
tokens = [lemmatizer.lemmatize(word) for word in tokens]
# pos-tagging
pos_tags = nltk.pos_tag(tokens)

nouns = ""
for i in pos_tags:
    if i[1] == 'NN' or i[1] == 'NNP' or i[1] == 'NNS' or i[1] == 'NNPS':
        nouns = nouns + " " + i[0]

wc = WordCloud(max_words=20).generate(nouns)
plt.imshow(wc, interpolation='bilinear')
plt.axis("off")
plt.show()
```
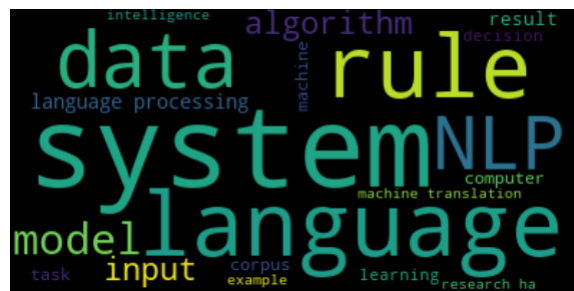


Figura 1: Word Cloud da questão 1.

## 1.2 Gerar um word cloud contendo os 20 lemas de verbos mais frequentes.

```
verbs = " ";
for i in pos_tags:
    if i[1] == 'VB':
        verbs = verbs + " " + i[0]

wc = WordCloud(max_words=20).generate(verbs)
```

```
7
8  plt.imshow(wc, interpolation='bilinear')
9  plt.axis("off")
10 plt.show()
```
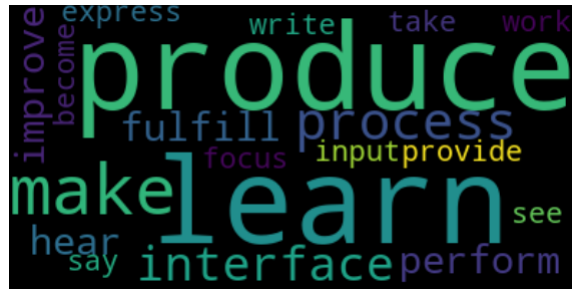


Figura 2: Word Cloud da questão 2.

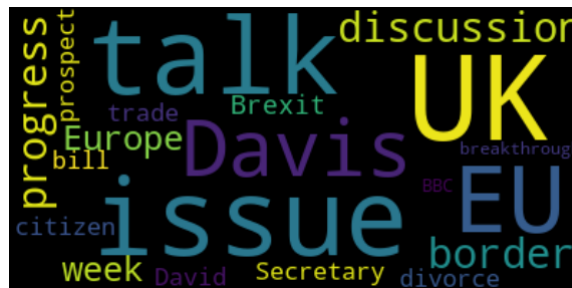## 1.3 Refazer os itens 1 e 2 usando a lista de stop words.


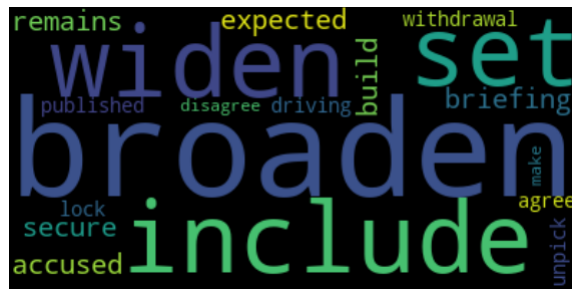
Figura 3: Word Cloud de substantivos da questão 3.



Figura 4: Word Cloud de verbos da questão 3.

```
1  arq = open('Corpus_en_NER.txt','r')
2  text = arq.read()
3  arq.close()
4
5  lemmatizer = WordNetLemmatizer()
6  # tokenizing
7  tokens = nltk.word_tokenize(text)
8  print(tokens)
9  # lemmatizing
10 tokens = [lemmatizer.lemmatize(word) for word in tokens]
11 print(tokens)
12 # pos-tagging
13 pos_tags = nltk.pos_tag(tokens)
14 sw = open('stopwords.txt','r')
```

```python
15  stopwords = sw.read()
16
17  nouns = ""
18
19  for i in pos_tags:
20      if i[1] == 'NN' or i[1] == 'NNP' or i[1] == 'NNS' or i[1] == 'NNPS':
21          if i[0] not in stopwords:
22              nouns = nouns + " " + i[0]
23
24  wc = WordCloud(max_words=20).generate(nouns)
25  plt.imshow(wc, interpolation='bilinear')
26  plt.axis("off")
27  plt.show()
28
29  verbs = "";
30  for i in pos_tags:
31      if i[1] == 'VB' or i[1] == 'VBD' or i[1] == 'VBG' or i[1] == 'VBN' or i[1] == 'VBP'
        or i[1] == 'VBZ':
32          if i[0] not in stopwords:
33              verbs = verbs + " " + i[0]
34
35  wc = WordCloud(max_words=20).generate(verbs)
36  plt.imshow(wc, interpolation='bilinear')
37  plt.axis("off")
38  plt.show()
```

## 1.4 Gerar uma árvore de parsing da frase.

```python
1  from nltk import Tree
2  from pycorenlp import StanfordCoreNLP
3
4  nlp = StanfordCoreNLP('http://localhost:9000')
5
6  text = "The last love letter I wrote was probably about 10 years ago."
7
8  output = nlp.annotate(text, properties={
9      'annotators': 'parse',
10     'outputFormat': 'json'
11 })
12
13 tree1 = output['sentences'][0]['parse'] + ""
14 tree2 = Tree.fromstring(tree1)
15 tree2.draw()
```
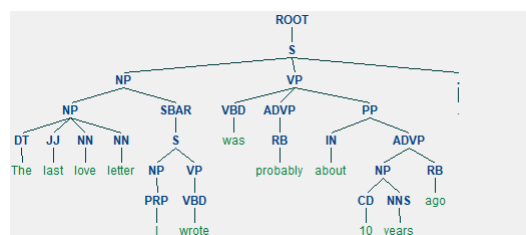


Figura 5: Parse tree da questão 4.

# 2 Usando o CoreNLP

```python
1  from stanfordcorenlp import StanfordCoreNLP
2  import json
3
4  nlp = StanfordCoreNLP("C:\ stanford−corenlp−full−2018−02−27")
5  arq = open("Corpus_en_NER.txt", "r")
```

```
 6  sentence = arq.read()
 7
 8  print("——— PIPELINE ———")
 9  print('[INFO] Tokenize:', nlp.word_tokenize(sentence))
10
11  props = {'annotators': 'tokenize,ssplit', 'pipelineLanguage': 'en', 'outputFormat': '
        json'}
12  sentence_splitting = nlp.annotate(sentence, properties=props)
13  jsonToPython = json.loads(sentence_splitting)
14  sentences = []
15  for sent in jsonToPython["sentences"]:
16      start_offset = sent['tokens'][0]['characterOffsetBegin']
17      end_offset = sent['tokens'][-1]['characterOffsetEnd']
18      sent_str = sentence[start_offset:end_offset]
19      sentences.append(sent_str)
20  print('[INFO] Sentence Splitting:', sentences)
21
22  pos = nlp.pos_tag(sentence)
23  print('[INFO] Part of Speech:', pos)
24
25  props = {'annotators': 'tokenize,lemma', 'pipelineLanguage': 'en', 'outputFormat': 'json
        '}
26  lemma = nlp.annotate(sentence, properties=props)
27  jsonToPython = json.loads(lemma)
28  lemmas = {}
29  for sent in jsonToPython["sentences"]:
30      for token in sent['tokens']:
31          word_original = token['word']
32          lemma_original = token['lemma']
33          if word_original not in lemmas:
34              lemmas[word_original] = lemma_original
35  print('[INFO] Lemmatization:', lemmas)
36
37  entities = nlp.ner(sentence)
38  print('[INFO] Named Entities:', entities)
39  print('[INFO] Dependency Parsing:', nlp.dependency_parse(sentence))
```

## 2.1 Os verbos identificados e os seus respectivos lemmas

```
 1  question1 = []
 2  for world in pos:
 3      if 'VB' in world[1]:
 4          world_aux = list(world)
 5          world_aux.append(lemmas[world[0]])
 6          question1.append(world_aux)
 7  print("Verbos e seus lemas: ", question1)
 8
 9  past_form = []
10  for world in pos:
11      if 'VBD' in world[1] or 'VBN' in world[1]:
12          world_aux = list(world)
13          world_aux.append(lemmas[world[0]])
14          past_form.append(world_aux)
15  print("Verbos no passado e seus lemas: ", past_form)
16  print("O lemma    o verbo no presente.")
```

## 2.2 As entidades nomeadas encontradas

```
 1  question2 = []
 2  for entiti in entities:
 3      if 'O' not in entiti[1]:
 4          question2.append(entiti)
 5  print(question2)
 6  arq.close()
```

## 2.3 Os tipos distintos de dependências gramaticais encontradas no arquivo NLP.txt

```
1  arq = open("NLP.txt", "r")
2  sentence = arq.read()
3  dependency = nlp.dependency_parse(sentence)
4  types = []
5  for dependences in dependency:
6      if dependences[0] not in types:
7          types.append(dependences[0])
8  print(types)
9  arq.close()
```

# 3  Implementando um sumarizador simples de notícias

```
1  import nltk
2  from glob import glob
3
4  def pipeline(text):
5      sent_detector = nltk.data.load('tokenizers/punkt/english.pickle')
6      sentences = sent_detector.tokenize(text.strip())
7      qtt = []
8      for sentence in sentences:
9          tokens = nltk.word_tokenize(sentence)
10         postag = nltk.tag.pos_tag(tokens)
11         parse_tree = nltk.ne_chunk(nltk.tag.pos_tag(tokens), binary=True)  # POS tagging
        before chunking!
12
13         named_entities = []
14
15         soma = []
16
17         for t in parse_tree.subtrees():
18             if t.label() == 'NE':
19                 named_entities.append(list(t))  # if you want to save a list of tagged
        words instead of a tree
20                 soma.append(list(t).__len__())
21
22
23         named_entities.clear()
24         qtt.append(sum(soma))
25     return sentences, qtt
26
27
28  def getScore(sentences, qtt):
29      scoreSi = []
30      i = 0
31      n = qtt.__len__()
32      while i < qtt.__len__():
33          score = 1- ((i+1)/n)
34          scoreGlobal = 1 + ((2*qtt[i])/(n+score))
35          scoreSi.append([sentences[i], scoreGlobal])
36          i = i+1
37      return scoreSi
38
39
40  def summarized(score):
41      score.sort(key=lambda x: x[1], reverse = True)
42      n = score.__len__()
43      k = 0.3*n
44      k = int(float(k))
45      for i in range(0,n-k):
46          score.pop()
47
48      summarize = ""
49      for sentence in score:
50          summarize = summarize + sentence[0]
51
52      return summarize
53
54
55
```

```
56  def main():
57      for filepath in glob('News\\**'):
58          arq = open(filepath, 'r')
59          text = arq.read()
60          (sentences, qtt) = pipeline(text)
61          score = getScore(sentences, qtt)
62          print(summarized(score), "\n")
63
64  if __name__ == "__main__":
65      main()
```

## 3.1 A sumarização é razoável? Apresentam problema de coerência? Qual sua opnião?

O resumo trás, de fato, uma certa parcela das informações mais importantes no texto. Mas, apesar disso, sua informação fica muito dispersa pela falta de coesão entre as sentenças. As informações parecem jogadas, bagunçadas e não se complementam. De fato, resume razoavelmente bem o que o texto possui, porém, não resume bem o que ele quer dizer.