

1ª Lista de Exercícios - Mineração de Texto

Maria Fernanda Souza Andrade - fsandrade25@gmail.com

Saulo Lucas Gomes Ferreira - saulolucas.gm@gmail.com

2 de outubro de 2018

1 Quais as vantagens e desvantagens dos tipos de anotação de corpora inline vs. stand-off?

Uma das vantagens observadas na anotação *inline* comparada a anotação *stand-off*, é o metadado já está acompanhado no mesmo arquivo que o dado, assim a extração dessa informação torna-se mais rápida comparando com a outra anotação. No entanto, essa forma de estruturação também se mostra como uma desvantagem, pois é necessário realizar um pré-processamento com o intuito de separar o dado do metadado.

Em contra partida para a anotação *stand-off* a sua vantagem é exatamente o seus metadados já estarem separados em um arquivo diferente, entretanto, a necessidade de realizar referências entre esses dois arquivos torna a leitura custosa aparecendo como uma desvantagem para o método.

2 Ao seu ver, porque o formato híbrido, isto é, a anotação "em camadas" é o mais usado hoje em dia?

A organização sistemática feita dessa forma evita que se suje ou acumule informações demais no conteúdo principal. Separar dados dá mais legibilidade e, consequentemente, diminui os erros. Utilizando índices para ligar diferentes camadas mantém os dados divididos de acordo com a sua função na anotação e evita mistura de informações.

3 Quais tipos de anotação são realizadas pela ferramenta BRAT (<http://brat.nlplab.org>)? Para quais problemas de mineração de textos ela pode ser usada?

Essa ferramenta possibilita realizar marcações nas palavras presentes no texto informando o tipo de entidade que elas estão vinculadas, além disso, ela também realiza uma análise de dependência entre palavras(anotações) de uma sentença, normaliza palavras, permite classificar as frases como sendo substantiva, verbal ou preposicional. Sendo assim, essa ferramenta é mais utilizada para problemas do tipo de extração de informação e processamento de linguagem natural.

4 O que se entende por POS tagging? E qual a diferença entre POS taging e parsing?

POS tagging contém funções gramaticais simples que são utilizadas para marcar separadamente cada palavra que compõe um texto, enquanto que o *parsing* se baseia em dependências binárias entre palavras.

5 Discuta sobre as propriedades do parsing constituinte e o parsing de dependências.

São duas diferentes maneiras de se obter a constituição sintática de um texto, que se baseiam em diferentes conexões de papéis semânticos de uma palavra num determinado contexto.

O *parsing constituinte* faz uma divisão hierárquica do contexto geral até a separação atômica das palavras e a determinação do seu papel na frase. É, visualmente, mais prático e fácil de se analisar humanamente, tendo em vista que a pirâmide que se é obtida na separação gradativa é bem explícita nas camadas que separam determinadas classes. No entanto, se mostra incapaz de relacionar semanticamente duas palavras atômicas de forma direta, o faz em contexto de função de dois campos semânticos no

período, como a relação de sujeito e predicado, mas a relação direta entre palavras, na maior parte das vezes, fica implícita.

Se a explicitação dessa relação for necessária, o *parsing de dependências* se mostra mais eficaz para essa tarefa, pois tem foco em como as palavras se relacionam. Esse modelo conecta palavras com dependências semânticas numa oração. Todavia, esse modelo parece complexo ocasionalmente até para humanos, quem dirá máquinas.

6 Um algoritmo de tokenização para a língua portuguesa.

```
1 public static String getNextToken() {
2     char ch;
3     String token = new String();
4
5     while(true){
6         cursor = currentPosition;
7         if(stream.length() >= cursor) ch = stream.charAt(cursor);
8         else return null;
9
10        while(cursor < stream.length()-1 && !isDelimiterSet(ch)){
11            ch = stream.charAt(++cursor);
12        }
13
14        if(cursor >= stream.length()){
15            cursor = currentPosition;
16            return null;
17        }
18
19        if(isEmptySpace(ch)){
20            if(cursor == currentPosition){
21                currentPosition+=1;
22                return null;
23            }
24            else{
25                token = stream.substring(currentPosition, cursor);
26                currentPosition = cursor+1;
27                return token;
28            }
29        }
30
31        if(cursor == currentPosition){
32            token = ch + " ";
33            currentPosition = cursor+1;
34        } else {
35            token = stream.substring(currentPosition, cursor);
36            currentPosition = cursor;
37            return token;
38        }
39
40        return null;
41    }
42 }
43
44 public static boolean isDelimiterSet(char ch){
45     if(delimiterSet.contains(ch + " ")) return true;
46     return false;
47 }
48
49 public static boolean isEmptySpace(char ch){
50     if(whiteSpace.contains(ch + " ")) return true;
51     return false;
52 }
```

7 Um algoritmo para segmentação de sentenças(orções) de textos em português.

```
1 public static String segmentation() {
2     StringBuilder text = new StringBuilder(stream);
3
4     char ch;
5     for(int i = 0; i < text.length()-1; i++){
6         ch = text.charAt(i);
7         if(text.length() == '?' || text.length() == '!'){
8             text.replace(i+1, i+2, "\n");
9         }
10    }
```

```

10         else if((ch == '\\' || ch == '\n') && text.charAt(i+1) == '.'){
11             text.replace(i+1, i+2, "\\n");
12
13         }
14         else if(text.charAt(i+1) != ' '){
15
16         }
17         else if(ch == ']' || ch == ')' || ch == '}'){
18             if(text.charAt(i+1) == '.'){
19                 text.replace(i+1, i+2, "\\n");
20             }
21         }
22         else if(ch == '.'){
23             if(Character.isUpperCase(thisToken(i-1).charAt(0)) && thisToken(i-1).
length() < 5 && Character.isUpperCase(text.charAt(i+2))){
24
25             }else if(Character.isLowerCase(thisToken(i-1).charAt(0)) && Character.
isUpperCase(text.charAt(i+2))){
26                 text.replace(i+1, i+2, "\\n");
27             }else if(thisToken(i-1).length() < 2){
28
29             }else{
30                 text.replace(i+1, i+2, "\\n");
31             }
32         }
33         else if(ch == ' '){
34             if(text.charAt(i+1) == '(' || text.charAt(i+1) == '{' || text.charAt(i+1)
== '[' || text.charAt(i+1) == '\\ || text.charAt(i+1) == '"' || text.charAt(i+1)
== '$'){
35                 text.replace(i, i+1, "\\n");
36             }
37         }
38     }
39     return text.toString();
40 }
41
42 // Funcao que retorna o token daquela determinada posicao
43 public static String thisToken(int i){
44     while(!whiteSpace.contains(stream.charAt(i) + "") && !delimiterSet.equals(stream.
charAt(i) + "")){
45
46         i--;
47
48         //System.out.println("a");
49     }
50     int first = ++i;
51     while(!whiteSpace.contains(stream.charAt(i) + "") && !delimiterSet.equals(stream.
charAt(i) + "")){
52
53         i++;
54         //System.out.println("a");
55     }
56     return stream.substring(first, i);
57 }

```

8 Use o NLTK para criar um pipeline.

```

1  arq = open("texto-en.txt", "r")
2
3  texto = ""
4  texto = arq.read()
5
6  # ----- Tokenization -----
7  tokens = nltk.word_tokenize(texto)
8  print("[INFO] Tokenization:")
9  print(tokens)
10
11 # ----- Sentence Splitting -----
12 sent_detector = nltk.data.load('tokenizers/punkt/english.pickle')
13 sentences = sent_detector.tokenize(texto.strip())
14 print("[INFO] Sentence Splitting:")
15 print(sentences)
16
17 # ----- Lemmatization (Reduce inflections or variant forms to base form)-----
18 from nltk.stem import WordNetLemmatizer
19 wordnet_lemmatizer = WordNetLemmatizer()
20 lema = {}

```

```

21 for w in tokens:
22     if w not in lema:
23         lema[w] = wordnet_lemmatizer.lemmatize(w)
24 print("[INFO] Lematization:")
25 print(lema)
26
27 # ----- Stemming (Reduce terms to their stems or roots) -----
28 from nltk.stem.porter import PorterStemmer
29 porter_stemmer = PorterStemmer()
30 stem = {}
31 for w in tokens:
32     if w != ',' and w != '.' and w != '!' and w != '?' and w != '"' and w != "'" and w
    != '':
33         if w not in stem:
34             stem[w] = porter_stemmer.stem(w)
35 print("[INFO] Stemming:")
36 print(stem)
37
38 # ----- POS tagging -----
39 tagged = nltk.pos_tag(tokens)
40 print("[INFO] POS tagging:")
41 print(tagged)
42
43 arq.close()

1 # Respostas:
2 print("\n")
3 split = []
4 for w in tokens:
5     if w != ',' and w != '.' and w != '!' and w != '?' and w != '"' and w != "'" and w
    != '':
6         split.append(w)
7 print("Letra A: ", len(split)) # Resposta 693
8 print("Letra B: ", len(stem.keys())) # Resposta 326
9
10 tokens_per_sentence = []
11 for s in sentences:
12     token = nltk.word_tokenize(s)
13     tokens_per_sentence.append(len(token))
14
15 print("Letra C: {0} sentencas e {1:2.4f} media de tokens por sentenca".format(len(
    sentences), (sum(tokens_per_sentence)/len(tokens_per_sentence))))
16
17 print("Letra E: Grafico em anexo")
18
19 tag_fd = nltk.FreqDist(tag for (word, tag) in tagged)
20
21 freq = []
22 tags = []
23 for i in tag_fd.most_common():
24     tags.append(i[0])
25     freq.append(i[1])
26
27 # ----- Grafico questao 08 letra e -----
28 plt.figure(0)
29 plt.bar(tags, freq)
30 plt.xlabel("TAGS")
31 plt.ylabel("Frequency")
32 plt.title("Question 8 Letter E")
33 plt.grid(True)
34 fig = plt.gcf()
35 fig.set_size_inches(18.5, 10.5, forward=True)
36 plt.savefig('letter_e.png')
37
38 print("Letra F: ")
39 stem_cont = {}
40 for w in tokens:
41     if w != ',' and w != '.' and w != '!' and w != '?' and w != '"' and w != "'" and w
    != '':
42         w_aux = porter_stemmer.stem(w)
43         if w_aux not in stem_cont:
44             stem_cont[w_aux] = 1
45         elif w_aux in stem_cont:
46             cont = stem_cont[w_aux]
47             stem_cont[w_aux] = cont + 1
48
49 from operator import itemgetter
50 stem_cont = dict(reversed(sorted(stem_cont.items(), key=itemgetter(1))))

```

```

51 freq = list(stem_cont.values())
52 stems = list(stem_cont.keys())
53 print(freq)
54 print(stems)
55
56 # ----- Grafico questao 08 letra F -----
57 plt.figure(1)
58 plt.bar(stems[0:15], freq[0:15])
59 plt.xlabel("Steaming")
60 plt.ylabel("Frequency")
61 #plt.yticks(sorted(set(list(stem_cont.values()))))
62 plt.title("Question 8 Letter F")
63 plt.grid(True)
64 fig = plt.gcf()
65 fig.set_size_inches(18.5, 10.5, forward=True)
66 plt.savefig('letter_f.png')
67
68

```

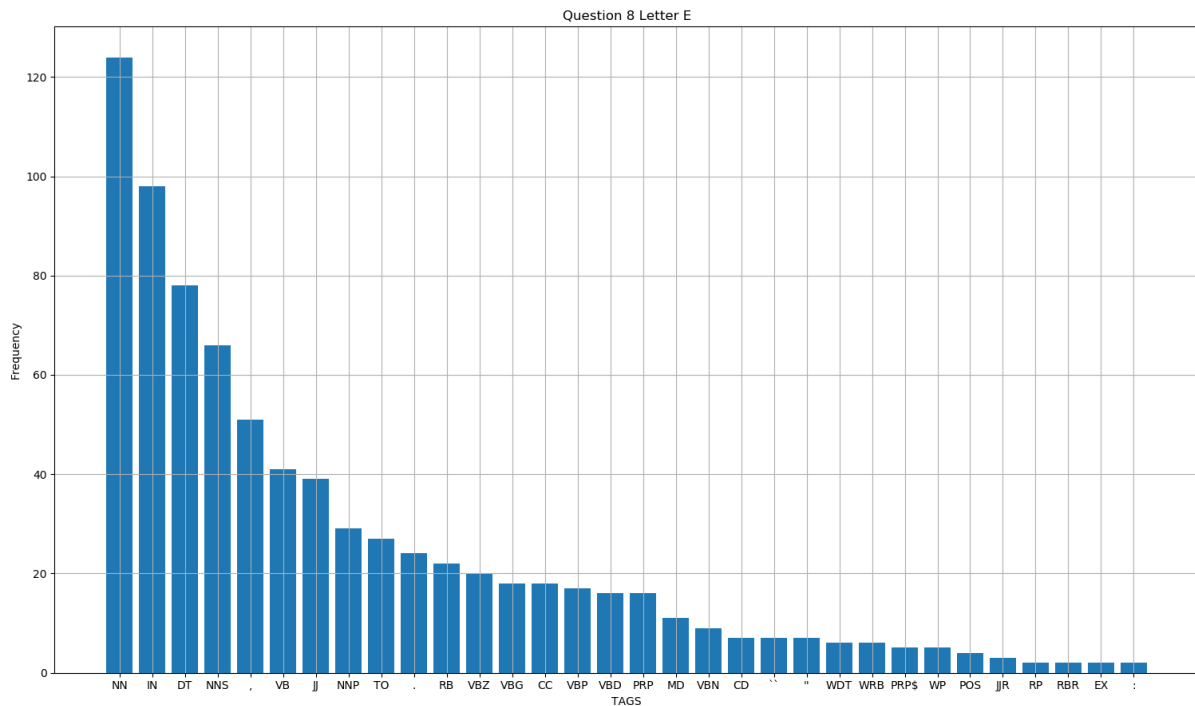


Figura 1: Gráfico referente a letra E da questão 8.

e) As classes que tem maior ocorrência são os substantivos, preposições e determinantes ou artigos.

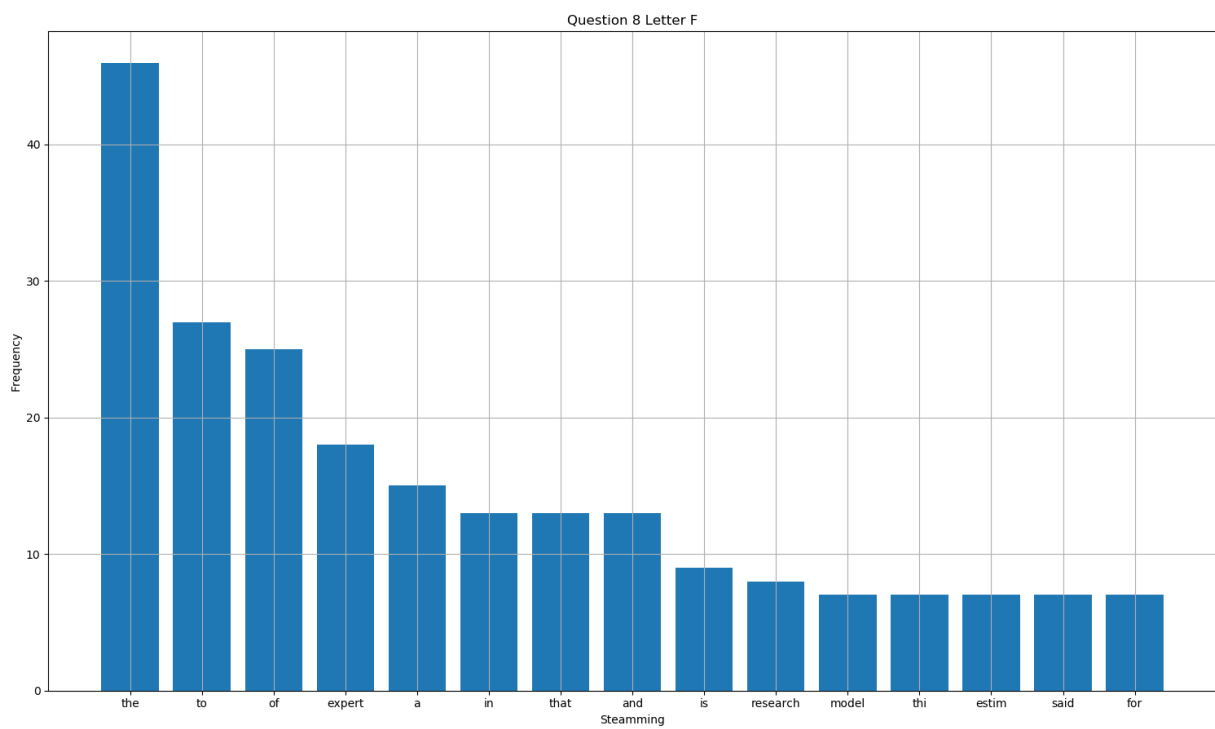


Figura 2: Gráfico referente a letra F da questão 8.

f) Os radicais mais frequentes são geralmente artigos e preposições. Por se tratarem de palavras padrões de ligação, é comum que se use com mais frequência no texto e, conseqüentemente, que se identifique mais.