

# SEMAFOR'S DATA STATISTICAL PARSING DOCUMENTATION

## Summary

1. Context
2. Tools
3. Python Modules
4. Code
5. Results
6. Conclusion

## Context

This parsing starts from a point where we have semafor's output as a [Json](#) file for each abstract of our corpus, this Json have the following structure: each line has a Json object with the attribute "frames", this attribute have an [Array](#) that brings us information about the identified frames, each position in the Array is a identified frame, this information comes in two attributes: "target" and "annotationSets".

"target" is where general information about the frame is, his name and also the name and position (counting spaces) of the lexical unit that called it, in the attributes "start", "end" and "text".

"annotationSets" is where we will find information about the result of semafor's analysis, like "rank", "score", that is the score that semafor gave to the frame\*, and "frameElements" that are the frame elements identified for that frame.

The first line in this Json files brings the frames identified in the name of the journal, the ones identified for the name of the abstract will be find at the second line, the rest of the file have the ones identified for the abstract itself.

The objective of this parsing is to extract quantitative and qualitative statistical analysis for this semafor results, initially, these are:

1. Occurrence of frames in the corpus, for all abstracts, per field, per discipline and per journal
2. Occurrence of Lexical Units in the corpus, for all abstracts, per field, per discipline, per journal and also for each frame (considering frame subdivisions)

## Tools

The tools used to extract this analysis were the programming language [Python](#) and the [Jupyter Notebook](#) environment. Python is well known for having great performance in processing large amounts of data, this will be helpful in dealing with our corpus, Jupyter allow us to have a practical environment to generate our statistical info, with the concept of live code that can be used to run code from a point where data have already been treated, without needing to run all the code every time.

## Python Modules

The following Python libraries were used to facilitate some tasks in the code:

Os:

This module is used to be possible to pass some commands to the operational system, in our context it was used to handle directories

Sys:

This module is used to access some operational system variables, in our context it was used also to handle directories

[Matplotlib](#):

This module can be used to many things (check documentation linked), we used it to deal with charts plotting.

Json:

This module is used to handle Json formatted text and objects, we used it to handle the Json objects that Semafor outputs.

Math:

This module handles math calculations and we used it exactly for that.

[Inflection](#):

This module can be used to deal with linguistics problems in text processing, we used it to get the singular of words.

## Code

Here are some of the most important parts of the code.

### itemsPerCategory function:

This function can receive a minimum of two e maximum of four parameters, the first two are an array with the file paths of semafor output and the item type that it will search for (frames or lexical units), the other two are optional parameters that, if used, will be an array with the categories (names of journals, fields or disciplines) and the category type (by journal, field or discipline). Basically it will go through the files that match the categories using the **countItems** function to count the occurrence of each item, respecting the item type.

```
1 def itemsPerCategory(files, item_type, *params):
2     mensagem = "Dados de: \n" #log
3     result = dict() #output
4
5     #considering all files
6     if not params:
7         print("Searching for "+item_type+"s in all files")
8         all_result = dict()
9         #all_result = first30(sortDict(countItems(item_type,all_result)))
10        all_result = sortDict(countItems(item_type,all_result))
11        qtd = len(all_result)
12        cutstart = 0
13        while(cutstart < qtd):
14            generateChart(cutdict(all_result,cutstart,30), 'all', 'all ' + str(cutstart)+"_"+str(cutstart+30), item_type)
15            cutstart = cutstart+30
16
17        return all_result
18    else:
19        categories = params[0]
20        category_type = params[1]
21        print("Searching for "+item_type+"s per "+category_type+"s")
22        all_categories_result = dict()
23        for category in categories:
24            category_result = dict()
25            #category_result = first30(sortDict(countItems(item_type, category_result, category, category_type)))
26            category_result = sortDict(countItems(item_type, category_result, category, category_type))
27            if(category_result):
28                label = 'Lexical Units' if (item_type == 'lexical_unit') else 'Frames'
29                qtd = len(category_result)
30                cutstart = 0
31                while(cutstart < qtd):
32                    generateChart(cutdict(category_result,cutstart,30), category_type, category+" (" +label+" ) "+str(
33                        cutstart)+"_"+str(cutstart+30), item_type)
34                    cutstart = cutstart+30
35                all_categories_result[category] = category_result
36            else:
37                print("Nothing found for "+category_type+": "+category)
38        return all_categories_result
```

## countItems function:

```
1 def countItems(item_type, items_counted, "params"):
2     if params:
3         category = params[0]
4         category_type = params[1]
5         #traversing all files
6         for file_name in files:
7             file_location = path+"/"+file_name
8             abstract = getInfoFromFile(file_location)
9             file_content = open(file_location, "r").readlines()
10
11             #traversing all lines in each file (each line has a json with information about frames)
12             for line in file_content:
13                 jsonline = json.loads(line)
14
15                 #verifying each item
16                 for frame in jsonline.get('frames'):
17
18                     if(item_type == 'frame'):
19                         frame_name = frame.get('target').get('name').lower()
20                         if(frame_name[len(frame_name)-1] == 's' or frame_name[len(frame_name)-4:len(frame_name)-1] == 'ies'):
21                             frame_name = inflection.singularize(frame_name)
22                         if params:
23                             if(abstract.get(category_type) == category):
24                                 if(items_counted.get(frame_name) != None):
25                                     items_counted[frame_name] += 1
26                                 else:
27                                     items_counted[frame_name] = 1
28                             else:
29                                 if(items_counted.get(frame_name) != None):
30                                     items_counted[frame_name] += 1
31                                 else:
32                                     items_counted[frame_name] = 1
33
34                     if(item_type == 'lexical_unit'):
35                         for span in frame.get('target').get('spans'):
36                             lexical_unit = span.get('text').lower()
37                             if(lexical_unit[len(lexical_unit)-1] == 's' or lexical_unit[len(lexical_unit)-4:len(lexical_unit)-1] == 'ies'):
38                                 lexical_unit = inflection.singularize(lexical_unit)
39                             if params:
40                                 if(abstract.get(category_type) == category):
41                                     #se ja existir o index, soma
42                                     if(items_counted.get(lexical_unit) != None):
43                                         items_counted[lexical_unit] += 1
44                                     #senao cria o index sem modificacoes
45                                 else:
46                                     items_counted[lexical_unit] = 1
47                             else:
48                                 if(items_counted.get(lexical_unit) != None):
49                                     items_counted[lexical_unit] += 1
50                                 else:
51                                     items_counted[lexical_unit] = 1
52
53     return items_counted
```

## getTitle function:

This function is used to get the title of the abstract from the raw text file.

```
1 def getTitle(title):
2     txtFiles = list()
3     for (dirpath, dirnames, filenames) in os.walk("../raw"):
4         txtFiles += [os.path.join(dirpath, file) for file in filenames]
5
6     for txtFile in txtFiles:
7         if txtFile.endswith(".txt"):
8             for line in list(open(txtFile)):
9                 if(line.startswith("TI ")):
10                     nextline = list(open(txtFile))[list(open(txtFile)).index(line)+1]
11                     if(not nextline.startswith("SO ")):
12                         fileTitle = line.replace("TI ", "").replace("\n", " ").replace(" ", " ") + " ".join(nextline.split())
13                     else:
14                         fileTitle = line.replace("TI ", "").replace("\n", " ")
15                     if(fileTitle == title):
16                         return fileTitle
17
```

### getJournal and searchJournal functions:

This functions get the Journal of an abstract based on the file name

```
1 def searchJournal(file1,file2):
2     retorno = ""
3     for j in getAllJournals():
4         if(file1.find(j) != -1):
5             retorno = file1
6         elif(file2.find(j) != -1):
7             retorno = file2
8         elif( (file2.find(j) != -1) and (file1.find(j) != -1) ):
9             retorno = "False"
10    return retorno
11
12 def getJournal(file):
13     file = file.split('.seg.json')[0]
14     file1 = file.split('_')[-3].replace('-', ' ').upper()
15     file2 = file.split('_')[-2].replace('-', ' ').upper()
16     journal = searchJournal(file1,file2)
17     return journal
```

### getInfoFromFile function:

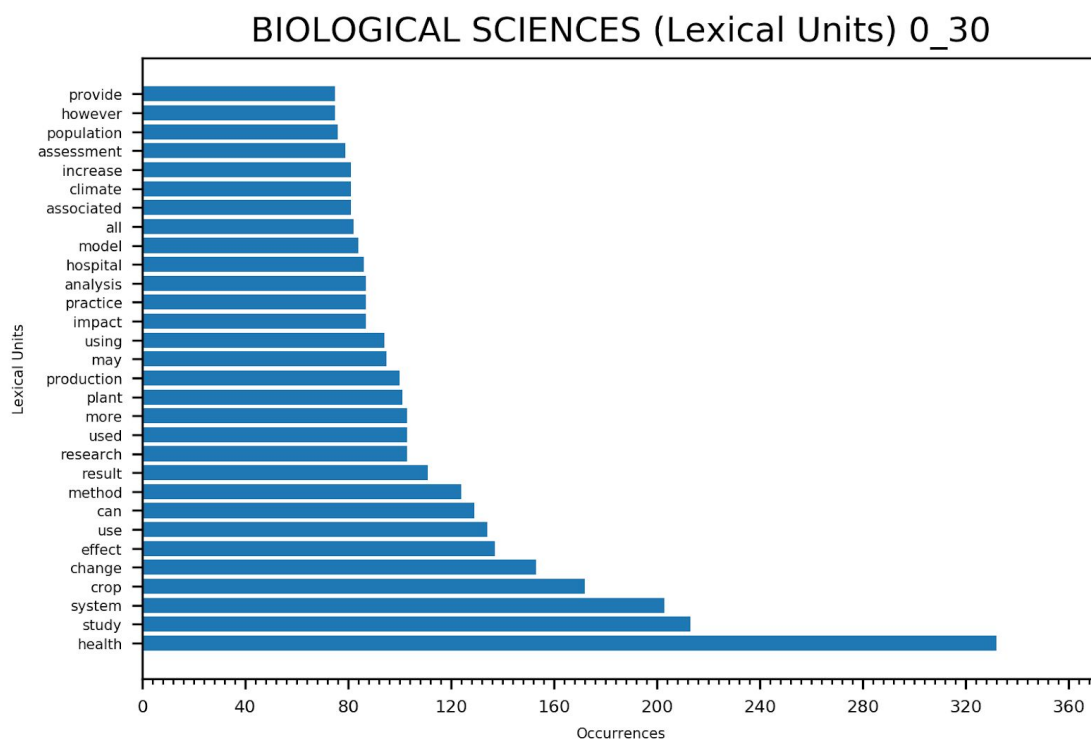
This function gets abstract information (journal, field, discipline) from a json file 'file'

```
1 def getInfoFromFile(file):
2     fields = dictionaries.dictionaryFields()
3     disciplines = dictionaries.dictionaryDisciplines()
4     journals = getAllJournals()
5     journal = ""
6
7     journal = getJournal(file)
8
9     field = fields.get(journal)
10    discipline = disciplines.get(field)
11    abstract = dict()
12    abstract['journal'] = journal
13    abstract['field'] = field
14    abstract['discipline'] = discipline
15    return abstract
```

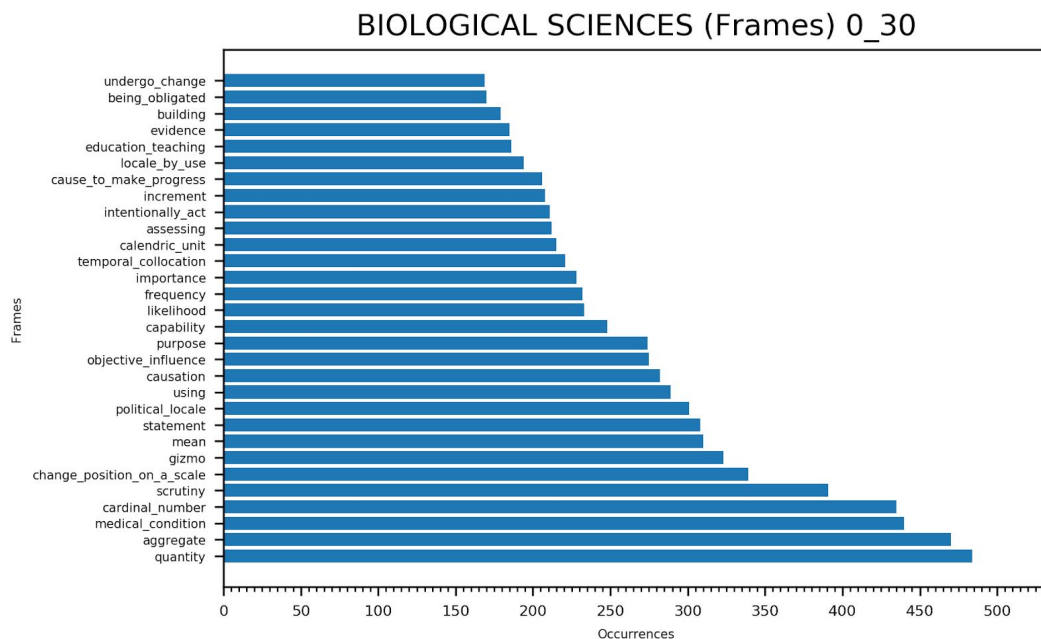
## Results

The objective of this work is to generate statistical analysis about semafor's output, the complete code and the results we already have will be available in [our GitHub repository](#), here's some examples of our results:

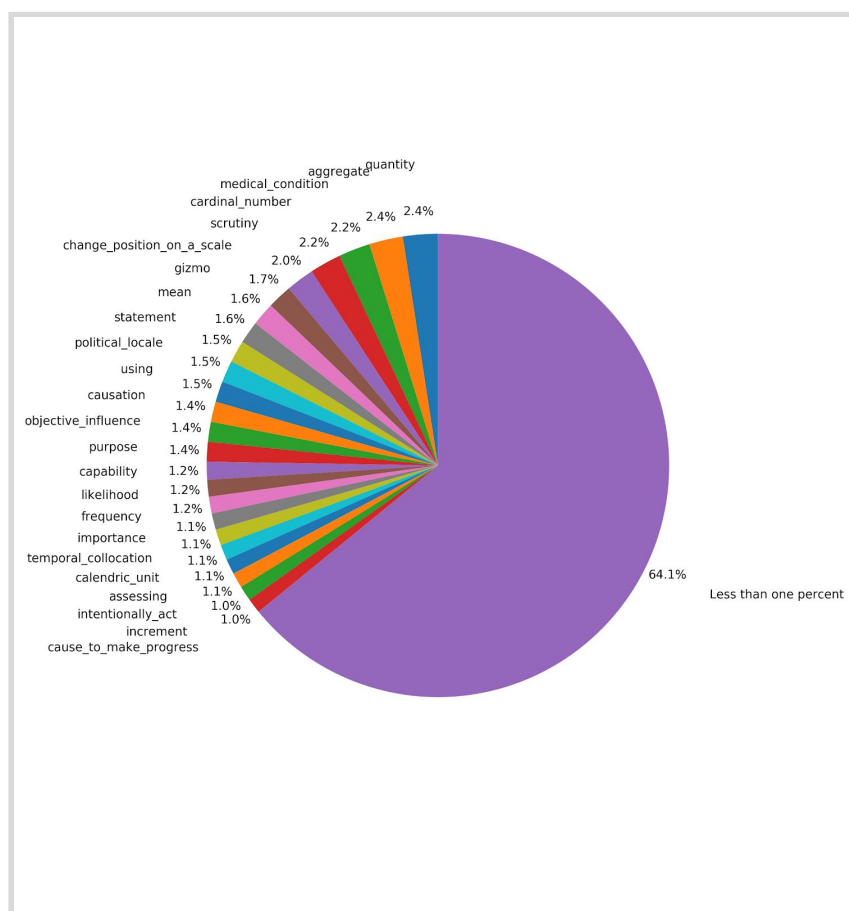
Counting of the thirty lexical units that most evoke frames in Biological Sciences (discipline):



Counting of the thirty frames that are more evoked in Biological Sciences (Discipline):



Percentage of the thirty more evoked frames in Biological Sciences (Discipline), grouping frames with less than one percent of occurrence:



Percentage of the more evoked frames in Biological Sciences (Discipline), considering only frames with more than one percent of occurrence:



