

The Problem

For this assignment, I want to create an algorithm that stores the words from the lyrics of your favorite songs. This is a very interesting algorithm, as it shows different applications of data structures, and helps me practice, what I believe, is one of the most important skills for coding, working with pre-existing libraries and API.

For this application, I'll be comparing 2 different data structures for storing words and counting them, and I will also work with the Tekore library, an API that connects python with data from Spotify, and the Lyrics Genius library, another API that connects python with the data from the website Genius, to get the lyrics of songs.

Modifying the trie structure

In this first part, I'm using the feedback from my previous assignment to modify the Trie data-structure for it to be more efficient.

The first modification was storing the children nodes in dictionaries (hash tables) instead of lists, which decreases the time complexity of the lookup from $O(n)$ to $O(1)$.

The second modification was using a heap (with the heapq library) to store the list of repetitions instead of a list. This decreases the time complexity of the k_most_common from $O(n \log n)$ to $O(k \log n)$. That is because by storing in a list we have to sort the whole list to get the k most common, and even the best sorting algorithm we've seen has a time complexity of $O(n \log n)$. However, if we store on a heap, we have $O(1)$ to get the smallest (or biggest) element and $O(\log n)$ to take this element out of the heap in a way that we keep the heap properties, so if we do this k times, we have $O(k \log n)$.

In [1]:

```

1 import heapq
2
3 class Node:
4     """This class represents one node of a trie tree.
5
6     Parameters
7     -----
8     The parameters for the Node class are not predetermined.
9     However, you will likely need to create one or more of them.
10    """
11
12    def __init__(self, data = None, parent = None, count = 0):
13
14        self.data = data
15        self.parent = parent
16        self.children = {}
17        self.word_end = False
18        self.word_count = count
19
20 class Trie:
21     """This class represents the entirety of a trie tree.
22
23     Parameters
24     -----
25     The parameters for Trie's __init__ are not predetermined.
26     However, you will likely need one or more of them.
27
28     Methods
29     -----
30     insert(self, word)
31         Inserts a word into the trie, creating nodes as required.
32     lookup(self, word)
33         Determines whether a given word is present in the trie.
34     insert_word_list(self)
35         inserts all the words in the word list in the trie.
36     list_of_repetitions(self, root, wordsInOrder=[]):
37         give a list of al the words in the trie and how many times they repeated
38     k_most_common(self, k):
39         Finds k words inserted into the trie most often.
40    """
41
42    def __init__(self, word_list = None):
43        """Creates the Trie instance, inserts initial words if provided.
44
45        Parameters
46        -----
47        word_list : list
48            List of strings to be inserted into the trie upon creation.
49        """
50        self.word_list = word_list
51        self.root = Node()
52
53        if word_list:
54            self.insert_word_list()
55
56    def insert(self, word):
57        """Inserts a word into the trie, creating missing nodes on the go.
58
59        Parameters
60        -----
61        word : str
62            The word to be inserted into the trie.
63        """
64        word = word.lower() #makes all the letters lower-case
65        current = self.root #current node is the root
66        for letter in word: #iterate through the word
67
68            #checking if the current letter exists in the current node children
69            if letter in current.children:
70                current = current.children[letter] #current update
71
72            else:
73                newLetter = Node(data=letter,parent=current) #create new node for the letter
74                current.children[letter] = newLetter #add the node to the children
75                current = newLetter #updates current
76
77        current.word_end = True #when its the last letter, set it to be the end of a word
78        current.word_count += 1 #increases the word count
79
80
81    def insert_word_list(self):
82        """inserts all the words in the word list in the trie.
83
84        Parameters
85        -----
86        None
87
88        Returns
89        -----

```

```

90     None
91     """
92     for word in self.word_list:
93         self.insert(word)
94
95
96     def lookup(self, word):
97         """Determines whether a given word is present in the trie.
98
99         Parameters
100        -----
101        word : str
102            The word to be looked-up in the trie.
103
104        Returns
105        -----
106        bool
107            True if the word is present in trie; False otherwise.
108        """
109
110        word = word.lower()    #makes all the letters lower-case
111        current = self.root    #current node is the root
112        for letter in word:    #iterate through the word
113
114            #checking if the current letter exists in the current node children
115            if letter in current.children:
116                current = current.children[letter] #current update
117
118            else:
119                return False
120
121        #checks if it is a prefix or a word
122        if current.word_end:
123            return True
124        else:
125            return False
126
127
128     def list_of_repetitions(self, root, wordsInOrder=[]):
129         """give a list of al the words in the trie and how many times they repeated
130
131         Parameters
132        -----
133        Node - root
134
135        List of tuples - wordsInOrder
136
137        Returns
138        -----
139        List of tuples
140        """
141
142        current = root    #current node is the root
143        if current.children != {}:    #base case (when current node has no children)
144            for child in list(current.children.values()):    #iterate through currents children
145                if child.word_end:    #checks if the child is the end of a word
146                    letter = child    #stores the child node in the var letter
147                    word = ''    #initializing the word
148                    while letter.parent is not None:    #goes up in the trie until it reaches the root
149                        word += letter.data    #appends the letter in the word
150                        letter = letter.parent    #updates the letter
151                    word = word[::-1]    #since we climbed up the trie, we need to invert the word
152                    heapq.heappush(wordsInOrder, (-child.word_count, word))    #inserts the word and the count in the
153                    self.list_of_repetitions(child, wordsInOrder)    # recursively calls the the function again using t
154        return wordsInOrder
155
156
157     def k_most_common(self, k):
158         """Finds k words inserted into the trie most often.
159
160         You will have to tweak some properties of your existing code,
161         so that it captures information about repeated insertion.
162
163         Parameters
164        -----
165        k : int
166            Number of most common words to be returned.
167
168        Returns
169        -----
170        list
171            List of tuples.
172
173            Each tuple entry consists of the word and its frequency.
174            The entries are sorted by frequency.
175
176        Example
177        -----
178        >>> print(trie.k_most_common(3))
179        [('the', 154), ('a', 122), ('i', 122)]

```

```
180
181     This means that the word 'the' has appeared 154 times in the inserted text.
182     The second and third most common words both appeared 122 times.
183     """
184     completelist = self.list_of_repetitions(self.root, wordsInOrder=[]) #get the whole priority queue of words
185     km = [heapq.heappop(completelist) for i in range(k)] #gets the k smallest (the values for repetition are negative)
186     km = [(i[1], -i[0]) for i in km]
187     return km
```

```
In [2]: ► 1 #simple testing
        2 wordbank = "oi oi a teste teste teste".split()
        3 trie = Trie(wordbank)
        4 trie.k_most_common(1)
```

```
Out[2]: [('teste', 3)]
```

Doing the same functions, but with a hash table instead of trie

As my second data-structure, I decided to go with a hash table. I'm simply using the built-in dictionaries since this is way more efficient and less time-consuming than trying to come up with my own hash functions and methods to solve collisions.

I think that using a hash table here is a great solution since it has $O(1)$ time complexity for lookup and for insertion. It is also better memory-wise since instead of having an object with a dictionary inside for every letter of a word, it simply has one key and one value for every different word.

```

In [3]: 1 class HT:
2         def __init__(self, word_list = None):
3             """Creates the Hash Table instance, inserts initial words if provided.
4
5             Parameters
6             -----
7             word_list : list
8                 List of strings to be inserted upon creation.
9             """
10            self.word_list = word_list
11            self.data = {}
12
13            if word_list:
14                self.insert_word_list()
15
16        def lookup(self, word):
17            """Determines whether a given word is present in the Hash Table.
18
19            Parameters
20            -----
21            word : str
22                The word to be looked-up.
23
24            Returns
25            -----
26            bool
27                True if the word is present; False otherwise.
28            """
29
30            word = word.lower()    #makes all the letters lower-case
31            return (word in self.data)
32
33
34        def insert(self, word):
35            """Inserts a word into the Hash Table.
36
37            Parameters
38            -----
39            word : str
40                The word to be inserted.
41            """
42
43            word = word.lower()    #makes all the letters lower-case
44
45            if self.lookup(word):
46                self.data[word]+=1    #adds to the counter of the word if the word is already there
47            else:
48                self.data[word]=1    #simply uses a new slot of the dictionary if the word is not already there
49
50
51        def insert_word_list(self):
52            """inserts all the words in the word list in the Hash Table.
53
54            Parameters
55            -----
56            None
57
58            Returns
59            -----
60            None
61            """
62            for word in self.word_list:
63                self.insert(word)
64
65        def k_most_common(self, k):
66            """Finds k words inserted into the Hash Table most often.
67
68            Parameters
69            -----
70            k : int
71                Number of most common words to be returned.
72
73            Returns
74            -----
75            list
76                List of tuples.
77
78                Each tuple entry consists of the word and its frequency.
79                The entries are sorted by frequency.
80            """
81
82            k_most_common_keys = heapq.nlargest(k, self.data, key=self.data.get)    #gets the k largest elements, by count
83            return [(i,self.data[i]) for i in k_most_common_keys]

```

```
In [4]: 1 #simple testing
2 wordbank = "oi oi a teste teste teste".split()
3 ht = HT(wordbank)
4 ht.k_most_common(1)
```

```
Out[4]: [('teste', 3)]
```

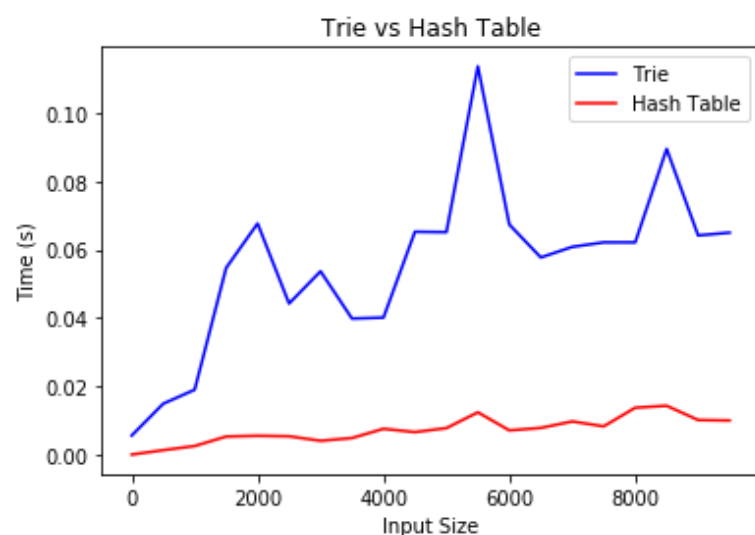
Comparing both data structures

In order to compare both structures, I'll use a more experimental approach. Using the big word bank from the previous assignment, I'll randomly (to make it less biased) select n words from it and measure the time that it takes for both structures to insert and give me the 1 most common and 1% (kind of) most common words, then I'll increase n and measure again. This will hopefully give me a graph that will model the asymptotic behavior of the functions.

```
In [5]: 1 #code adapted from the last assignment to get a good wordbank
2 import random
3 import string
4 import time
5 import matplotlib.pyplot as plt
6
7 import urllib.request
8 response = urllib.request.urlopen('http://bit.ly/CS110-Shakespeare')
9 bad_chars = [';', ',', '.', '?', '!', '1', '2', '3', '4', '5', '6', '7', '8', '9', '0', '_', '[', ']', '"', '&', "'"]
10
11 wordlist = []
12
13 for line in response:
14     line = line.decode(encoding = 'utf-8')
15     line = filter(lambda i: i not in bad_chars, line)
16     words = "".join(line).split()
17     for word in words:
18         wordlist.append(word)
```

```
In [6]: 1 #this cell might take a couple of minutes to run, please wait until it is done to run the next one
2
3 times_trie = []
4 times_HT = []
5
6 for i in range(1,10001,500):
7     wordbank = random.sample(wordlist, i)
8
9     time_trie = 0
10    time_HT = 0
11
12    for j in range(10):
13        s = time.perf_counter() #start timer
14        trie = Trie(wordbank)
15        trie.k_most_common(1)
16        trie.k_most_common(i//100)
17        e = time.perf_counter() #end timer
18        time_trie += (e-s)
19    times_trie.append(time_trie/10)
20
21    for j in range(10):
22        s = time.perf_counter() #start timer
23        ht = HT(wordbank)
24        ht.k_most_common(1)
25        ht.k_most_common(i//100)
26        e = time.perf_counter() #end timer
27        time_HT += (e-s)
28    times_HT.append(time_HT/10)
```

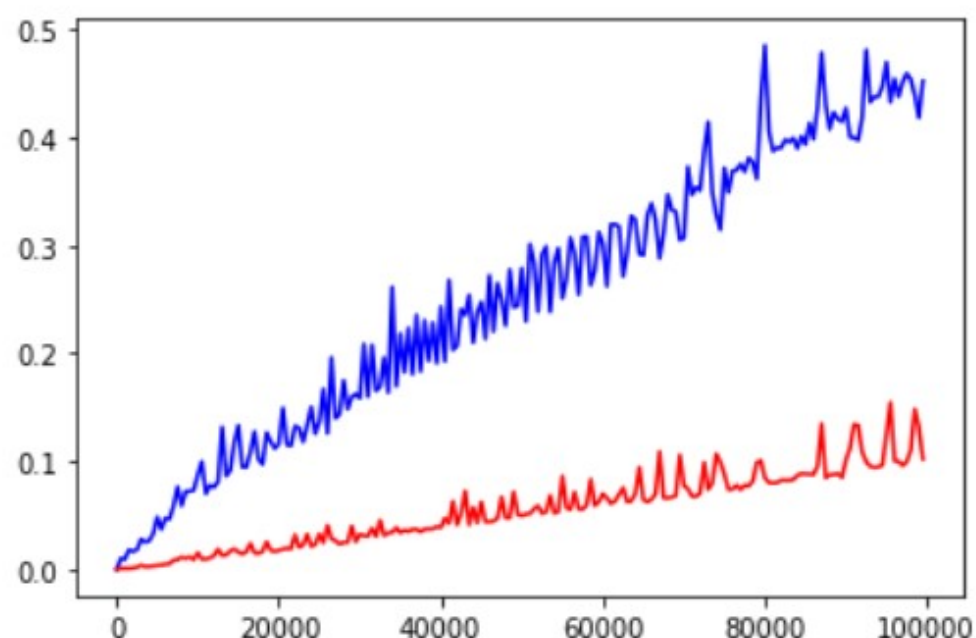
```
In [7]: 1 plt.plot(list(range(1,10001,500)),times_trie, c='blue')
2 plt.plot(list(range(1,10001,500)),times_HT, c='red')
3
4 plt.legend(['Trie','Hash Table'])
5 plt.xlabel('Input Size')
6 plt.ylabel('Time (s)')
7 plt.title('Trie vs Hash Table')
8
9 plt.show()
```



Since this code would take a very long time to run for a bigger number of words, here is an image of the same graph but with more data, so it will be easier to analyze.

```
In [12]: 1 from IPython.display import Image
2 Image(filename="graph.jpg")
```

Out[12]: Out[49]: [<matplotlib.lines.Line2D at 0x1de20258a58>]



The first we can see is that for every value of n , the hash table performs better. They both appear to have linear time complexity by this graph, however, the constant multiplier for the hash table is significantly lower.

If we analyze both more analitically, we know that they both have $O(k \log n)$ for the k_most_common , since they use the same basic "sort" to do it.

However, the hash table shines at the insertion by using a very simple constant complexity for each word, giving an overway complexiy of $O(n)$ for n words, while the Trie uses a more complicated code, which I meticulously described last assignment.

With these results, and the fact that the hash table code is simpler and more elegant, I chose to use it for my final implementation.

Experimenting with the Tekore API

I decide to use this open source library mainly for it's simplicity. It has a really complete documentation, well designed code that easilly connects you to the WEB API for Spotify , and even a Discord server to discuss problems.

<https://pypi.org/project/tekore/> (<https://pypi.org/project/tekore/>)

In [13]: ▶ 1 pip install tekore

Requirement already satisfied: tekore in c:\programdata\anaconda3\lib\site-packages (3.4.1)
 Requirement already satisfied: httpx<0.17,>=0.11 in c:\programdata\anaconda3\lib\site-packages (from tekore) (0.16.1)
 Requirement already satisfied: httpcore==0.12.* in c:\programdata\anaconda3\lib\site-packages (from httpx<0.17,>=0.11->tekore) (0.12.2)
 Requirement already satisfied: certifi in c:\programdata\anaconda3\lib\site-packages (from httpx<0.17,>=0.11->tekore) (2019.6.16)
 Requirement already satisfied: sniffio in c:\programdata\anaconda3\lib\site-packages (from httpx<0.17,>=0.11->tekore) (1.2.0)
 Requirement already satisfied: rfc3986[idna2008]<2,>=1.3 in c:\programdata\anaconda3\lib\site-packages (from httpx<0.17,>=0.11->tekore) (1.4.0)
 Requirement already satisfied: h11==0.* in c:\programdata\anaconda3\lib\site-packages (from httpcore==0.12.*->httpx<0.17,>=0.11->tekore) (0.11.0)
 Requirement already satisfied: idna; extra == "idna2008" in c:\programdata\anaconda3\lib\site-packages (from rfc3986[idna2008]<2,>=1.3->httpx<0.17,>=0.11->tekore) (2.8)
 Note: you may need to restart the kernel to use updated packages.

WARNING: You are using pip version 20.0.2; however, version 20.3.1 is available.
 You should consider upgrading via the 'C:\ProgramData\Anaconda3\python.exe -m pip install --upgrade pip' command.

In [14]: ▶ 1 import tekore as tk
 2
 3 #getting the authentication token
 4 conf = ("47721e5988094299b2a980889eea6403", "8bccab59ae3b428496a12572fcb36ad6", "https://sites.google.com/view/cs110-final-project/home?code=AQCe0J5aHKrOK3tIxtZEFmZNhydoAifY5Ha_cqeQymuu-VISMsJITXMxLcMIUh0RWxHA7-6zHrYRS09VZjKyAyKJD-b107v51aj73Gy2CmLP0W5nRNsIPqKQfuJfRGIPon-ESGq-bDeU3S5JWHkSeiT4xkerHUCYu8MJSOp1US55kVp_UbNq39EZSR_GeZdLI6Zo4r9Vs89Vp1QWQn0mx2tcW-1grqyTeyf8UuBms9vEfLvPMgTXTJGPML0qD-M7EndQJ5-4aQ12rRBVwckv35oTsM0FtRkY-1pGo7-iLaCPiRU1QHcCaI6Gj8P1EbidN86HkxK41UkbXOGQcd0qFoUoQaiq0cRIrwXp3klqvntwJA_RDzoj14xYY3A3qxjFFRf0BPNG_HDHypYlpV_wQiaIMBT7haLWg8r-SE7U9bs_mHe9Y6kXk1-jPRF5epB_HIttMFCeHQ7uJmCTZAF4x9TsYVMMU-83cZBifueRjq3fVvVibLGLvLiWkXEK3KsxJDAP4x5xpLWYoMGUV0cNIR9PKMomXsAHBjE1mP3HNJ5cQuFKB24maPRweWZuTXrnTunLx5KdA1QMmaPkuUEkluyNW53Y3Tidr_kkPlclDwVX4Z2Yk661nfKfLaSpusJ0QshkP8OC0m4kKsJpirPei94FJYyvMJtxRVHTOmXeITlDN7y0M6IpDQPh6UIPbr7jgCEtjKL-fevdVtCOQA&state=NgnZjkuykJJU_M1Doju4VdVD_W-G6g80uOXQVmCoBY (https://sites.google.com/view/cs110-final-project/home?code=AQCe0J5aHKrOK3tIxtZEFmZNhydoAifY5Ha_cqeQymuu-VISMsJITXMxLcMIUh0RWxHA7-6zHrYRS09VZjKyAyKJD-b107v51aj73Gy2CmLP0W5nRNsIPqKQfuJfRGIPon-ESGq-bDeU3S5JWHkSeiT4xkerHUCYu8MJSOp1US55kVp_UbNq39EZSR_GeZdLI6Zo4r9Vs89Vp1QWQn0mx2tcW-1grqyTeyf8UuBms9vEfLvPMgTXTJGPML0qD-M7EndQJ5-4aQ12rRBVwckv35oTsM0FtRkY-1pGo7-iLaCPiRU1QHcCaI6Gj8P1EbidN86HkxK41UkbXOGQcd0qFoUoQaiq0cRIrwXp3klqvntwJA_RDzoj14xYY3A3qxjFFRf0BPNG_HDHypYlpV_wQiaIMBT7haLWg8r-SE7U9bs_mHe9Y6kXk1-jPRF5epB_HIttMFCeHQ7uJmCTZAF4x9TsYVMMU-83cZBifueRjq3fVvVibLGLvLiWkXEK3KsxJDAP4x5xpLWYoMGUV0cNIR9PKMomXsAHBjE1mP3HNJ5cQuFKB24maPRweWZuTXrnTunLx5KdA1QMmaPkuUEkluyNW53Y3Tidr_kkPlclDwVX4Z2Yk661nfKfLaSpusJ0QshkP8OC0m4kKsJpirPei94FJYyvMJtxRVHTOmXeITlDN7y0M6IpDQPh6UIPbr7jgCEtjKL-fevdVtCOQA&state=NgnZjkuykJJU_M1Doju4VdVD_W-G6g80uOXQVmCoBY)
 5 token = tk.prompt_for_user_token(*conf, scope=tk.scope.every)
 6 spotify = tk.Spotify(token)
 7
 8 top_tracks = spotify.current_user_top_tracks(limit=10) #gets the 10 top tracks

Opening browser for Spotify login...

Please paste redirect URL: https://sites.google.com/view/cs110-final-project/home?code=AQCe0J5aHKrOK3tIxtZEFmZNhydoAifY5Ha_cqeQymuu-VISMsJITXMxLcMIUh0RWxHA7-6zHrYRS09VZjKyAyKJD-b107v51aj73Gy2CmLP0W5nRNsIPqKQfuJfRGIPon-ESGq-bDeU3S5JWHkSeiT4xkerHUCYu8MJSOp1US55kVp_UbNq39EZSR_GeZdLI6Zo4r9Vs89Vp1QWQn0mx2tcW-1grqyTeyf8UuBms9vEfLvPMgTXTJGPML0qD-M7EndQJ5-4aQ12rRBVwckv35oTsM0FtRkY-1pGo7-iLaCPiRU1QHcCaI6Gj8P1EbidN86HkxK41UkbXOGQcd0qFoUoQaiq0cRIrwXp3klqvntwJA_RDzoj14xYY3A3qxjFFRf0BPNG_HDHypYlpV_wQiaIMBT7haLWg8r-SE7U9bs_mHe9Y6kXk1-jPRF5epB_HIttMFCeHQ7uJmCTZAF4x9TsYVMMU-83cZBifueRjq3fVvVibLGLvLiWkXEK3KsxJDAP4x5xpLWYoMGUV0cNIR9PKMomXsAHBjE1mP3HNJ5cQuFKB24maPRweWZuTXrnTunLx5KdA1QMmaPkuUEkluyNW53Y3Tidr_kkPlclDwVX4Z2Yk661nfKfLaSpusJ0QshkP8OC0m4kKsJpirPei94FJYyvMJtxRVHTOmXeITlDN7y0M6IpDQPh6UIPbr7jgCEtjKL-fevdVtCOQA&state=NgnZjkuykJJU_M1Doju4VdVD_W-G6g80uOXQVmCoBY


In [15]: ▶ 1 #some songs have some info in the title that can make the website with the lyrics confused, like "Here I Go Again (2
 2 bad_for_song = ['(', '-']
 3
 4 for i in top_tracks.items:
 5 for j in range(len(i.name)-1):
 6 if i.name[j] in bad_for_song: #if I find a (or a -, I delete whatever comes after it from the song title
 7 i.name = i.name[:j-1]
 8 break
 9 print(i.name, [j.name for j in i.artists])

Relaxa! ['Haikaiss', 'Neo Beats', 'Cortesia da Casa', 'Ursso']
 Linda, Louca e Mimada ['Oriente', 'Rebeca']
 Pride and Joy ['Stevie Ray Vaughan']
 Sultans Of Swing ['Dire Straits']
 Céu Azul ['Charlie Brown Jr.']
 Philipa ['The Touré-Raichel Collective', 'Idan Raichel', 'Vieux Farka Touré']
 Here I Go Again ['Whitesnake']
 Beastly ['Vulfpeck']
 Zero ['Liniker e os Caramelows']
 Boredom ['Tyler, The Creator', 'Rex Orange County', 'Anna of the North']

Experimenting with LyricsGenius API


This was also the simplest API I found for working with lyrics of songs. Using an API of an existing database like this one is a much more effective and easier way of getting the job done without having to do a whole web search protocol from scratch in order to get some simple lyrics from a song.

<https://pypi.org/project/lyricsgenius/> (<https://pypi.org/project/lyricsgenius/>)

In [16]:  1 pip install lyricsgenius

```
Requirement already satisfied: lyricsgenius in c:\programdata\anaconda3\lib\site-packages (2.0.2)
Requirement already satisfied: beautifulsoup4>=4.6.0 in c:\programdata\anaconda3\lib\site-packages (from lyricsgenius) (4.7.1)
Requirement already satisfied: requests>=2.20.0 in c:\programdata\anaconda3\lib\site-packages (from lyricsgenius) (2.22.0)
Requirement already satisfied: soupsieve>=1.2 in c:\programdata\anaconda3\lib\site-packages (from beautifulsoup4>=4.6.0->lyricsgenius) (1.8)
Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in c:\programdata\anaconda3\lib\site-packages (from requests>=2.20.0->lyricsgenius) (1.24.2)
Requirement already satisfied: chardet<3.1.0,>=3.0.2 in c:\programdata\anaconda3\lib\site-packages (from requests>=2.20.0->lyricsgenius) (3.0.4)
Requirement already satisfied: certifi>=2017.4.17 in c:\programdata\anaconda3\lib\site-packages (from requests>=2.20.0->lyricsgenius) (2019.6.16)
Requirement already satisfied: idna<2.9,>=2.5 in c:\programdata\anaconda3\lib\site-packages (from requests>=2.20.0->lyricsgenius) (2.8)
Note: you may need to restart the kernel to use updated packages.
```

WARNING: You are using pip version 20.0.2; however, version 20.3.1 is available.
You should consider upgrading via the 'C:\ProgramData\Anaconda3\python.exe -m pip install --upgrade pip' command.

In [17]: 

```
1 import lyricsgenius
2
3 #token for the API
4 genius = lyricsgenius.Genius("CHIZ8mkPTjY4IJBKxF6dni13TDPlsBhxLN04YgXv1E13yBgGA800kpn3I35Eo0WK")
5
6 genius.response_format = 'plain' #setting the text as plain text
7
8 songs_with_lyrics = [genius.search_song(i.name, i.artists[0].name) for i in top_tracks.items] #gets the lyrics of m
9
10 songs_with_lyrics
```

```
Searching for "Relaxa!" by Haikaiss...
Done.
Searching for "Linda, Louca e Mimada" by Oriente...
Done.
Searching for "Pride and Joy" by Stevie Ray Vaughan...
Done.
Searching for "Sultans Of Swing" by Dire Straits...
Done.
Searching for "Céu Azul" by Charlie Brown Jr....
Done.
Searching for "Philipa" by The Touré-Raichel Collective...
No results found for: 'Philipa The Touré-Raichel Collective'
Searching for "Here I Go Again" by Whitesnake...
Done.
Searching for "Beastly" by Vulfpeck...
Done.
Searching for "Zero" by Liniker e os Caramelows...
Done.
Searching for "Boredom" by Tyler, The Creator...
Done.
```

Out[17]:

```
[('Relaxa!', 'Haikaiss'),
 ('Linda, Louca e Mimada', 'Oriente'),
 ('Pride and Joy', 'Stevie Ray Vaughan and Double Trouble'),
 ('Sultans of Swing', 'Dire Straits'),
 ('Céu Azul', 'Charlie Brown Jr.'),
 None,
 ('Here I Go Again', 'Whitesnake'),
 ('Beastly', 'Vulfpeck'),
 ('Zero', 'Liniker e os Caramelows'),
 ('Boredom', 'Tyler, The Creator')]
```

Combining all

In [18]:

```

1 class LyricsTable:
2     def __init__(self, word_list = None, bad_words = []):
3         """Creates the Hash Table instane, inserts initial words if provided.
4
5         Parameters
6         -----
7         word_list : list
8             List of strings to be inserted upon creation.
9
10        bad_words : list
11            List of words that won't be counted in the most common words.
12        """
13        self.word_list = word_list
14        self.data = {}
15        self.bad_words = bad_words
16
17        if word_list:
18            self.insert_word_list()
19
20    def lookup(self, word):
21        """Determines whether a given word is present in the Hash Table.
22
23        Parameters
24        -----
25        word : str
26            The word to be looked-up.
27
28        Returns
29        -----
30        bool
31            True if the word is present; False otherwise.
32        """
33
34        word = word.lower()    #makes all the letters lower-case
35        return (word in self.data)
36
37
38    def insert(self, word):
39        """Inserts a word into the Hash Table.
40
41        Parameters
42        -----
43        word : str
44            The word to be inserted.
45        """
46
47        word = word.lower()    #makes all the letters lower-case
48
49        if self.lookup(word):    #if the word is there
50            if word in self.bad_words:    #if it is one of the word I don't want to count in the list of repetitions
51                self.data[word] -= 1
52            else:
53                self.data[word] += 1
54        else:
55            if word in self.bad_words:    #if it is one of the word I don't want to count in the list of repetitions
56                self.data[word] = -1
57            else:
58                self.data[word] = 1
59
60
61    def insert_word_list(self):
62        """inserts all the words in the word list in the Hash Table.
63
64        Parameters
65        -----
66        None
67
68        Returns
69        -----
70        None
71        """
72        for word in self.word_list:
73            self.insert(word)
74
75    def k_most_common(self, k):
76        """Finds k words inserted into the Hash Table most often.
77
78        Parameters
79        -----
80        k : int
81            Number of most common words to be returned.
82
83        Returns
84        -----
85        list
86            List of tuples.
87
88            Each tuple entry consists of the word and its frequency.
89            The entries are sorted by frequency.

```

```

90         """
91
92         k_most_common_keys = heapq.nlargest(k, self.data, key=self.data.get) #gets the k largest elements, by count
93         return [(i,self.data[i]) for i in k_most_common_keys]
94
95     def repetitions(self, word):
96         """finds the amount of times a word shows up in the HT.
97
98         Parameters
99         -----
100        word : str
101            word you want to know the amount of times it shows up.
102
103        Returns
104        -----
105        int
106            number of times the word shows up.
107
108        """
109
110        if self.lookup(word): #if the word is in the HT
111            return abs(self.data[word]) #return the counter of the word (abs because the bad words are negative)
112        return 0
113
114     def initSpotify():
115         """Initializes the authentication for the Spotify API.
116
117         Parameters
118         -----
119         None
120
121        Returns
122        -----
123        None
124        """
125        conf = ("47721e5988094299b2a980889eea6403", "8bccab59ae3b428496a12572fcb36ad6", "https://sites.google.com/view/spotify-api")
126        token = tk.prompt_for_user_token(*conf, scope=tk.scope.everything)
127        spotify = tk.Spotify(token)
128
129     def getTopSongs(n, time_range):
130         """Gets your top tracks from Spotify.
131
132         Parameters
133         -----
134        n : int
135            number of tracks to get
136
137        time_range : str
138            range of the top track (long_term (forever), medium_term (last 6 months), short_term (last month))
139
140        Returns
141        -----
142        list
143            list of top tracks
144        """
145        top = spotify.current_user_top_tracks(limit=n, time_range=time_range) #gets top tracks
146
147        #some songs have some info in the title that can make the website with the lyrics confused, like "Here I Go Again"
148        bad_for_song = ['(', '-']
149
150        for i in top.items():
151            for j in range(len(i.name)-1):
152                if i.name[j] in bad_for_song: #if I find a ( or a -, I delete whatever comes after it from the song title
153                    i.name = i.name[:j-1]
154                break
155        return top
156
157     def getLyrics(song):
158         """Gets the lyrics of a song.
159
160         Parameters
161         -----
162        song : song object
163            song
164
165        Returns
166        -----
167        song object
168            song with the lyrics as an attribute
169        None
170            if no song
171        """
172        if song:
173            return genius.search_song(song.name, song.artists[0].name)
174        else:
175            return None
176
177     def List_of_lyrics(n, time_range):
178         """Creates a list of song lyrics from your top tracks on Spotify.
179

```

```

180     Parameters
181     -----
182     n : int
183         number you tracks to get
184
185     time_range : str
186         range of the top track (long_term (forever), medium_term (last 6 months), short_term (last month))
187
188     Returns
189     -----
190     list
191         list if top tracks with lyrics as an attribute
192     """
193
194     initSpotify()
195
196     top_songs = getTopSongs(n, time_range)
197
198     return [getLyrics(i) for i in top_songs.items]
199
200 def createWordbank(n, time_range):
201     """creates a bank of the words in the lyrics of your top tracks on Spotify
202
203     Parameters
204     -----
205     n : int
206         number you tracks to get
207
208     time_range : str
209         range of the top track (long_term (forever), medium_term (last 6 months), short_term (last month))
210
211     Returns
212     -----
213     list
214         list words
215     """
216     lyrics = List_of_lyrics(n, time_range) #list of songs with lyrics
217     wordbank = []
218     for i in lyrics: #iterates through the songs
219         if i: #if there is a song
220             #adds the words to the wordbank
221             wordbank += i.lyrics.replace("!", "").replace("?", "").replace(".", "").replace(",", "").replace(";", " ")
222     return wordbank

```

Feel free to play around with this part of the code and experiment with different numbers of top tracks, most common words, and time ranges. For any authentication errors, please contact me, it might be that the keys have expired and I have to provide new ones. If you do not use Spotify please refer to the output of the code in the pdf with outputs from my account.

In [19]:

1

wordbank = createWordbank(10,'long_term')

2

3

#excluding the words I don't want to count in my list of repetitions. Feel free to add more or take some out as you

4

bad_words = ['a','o','e','and','or','i','they','to','para','ele','ela','he','she','it','was','is','am','you','tu','

5

'pra','the','que','do','eu','de','é','of','she's',"he's",'me','um','se','te','that','on', 'in', "i'm",

6

'for', 'be','[chorus]','i've','so','this','if']

7

8

9

songTable = LyricsTable(wordbank,bad_words)

10

11

print(songTable.k_most_common(20))

Opening browser for Spotify login...
Please paste redirect URL: https://sites.google.com/view/cs110-final-project/home?code=AQAKDB1BLnDXqY05zh8WaB7W1ZhqXzm gTJsc85RfqoCSbvobdkPUBLIibSX2S1Th0VVgZY00yxUjEJRMXPel_pQegDpG4KQHxJUcDUB_DsbfnECL8qOVwis31FG4dxXbZyi5u0SnzPD42NF0cakyw xSibOgGaLdRT0d0LNLZQ3vq4n2ajLr-0YS2ZkC603rminveMx6dePloLI35494cdY7z6MPVJd5KU01_VeQey_QasEKoxTBPrU6FxOfpPCULMjx01E77LW5 FiV9VS3-WTmk89hcNINF8awVA12n8Gu8ynhBE3IUQUloXFEEIhnvufW27IH4y73twX3uo2DHW-eY7xPonq-WnDX76pZOV0iKQ_XhF_23ALTA7h0bEPJLig 2o_pds0BTSWFpyNQ-z9RuX-jq7y0vR_C6X0jUlF5mkB41nHbGq_N7zwDEEoNXCZ6CrjdLUdBgpMzcAZkuGuHanj95yYa2Tna6QFIa7VV2u7tic8kELrLCp 1A-NROD_PKviC4QLza81HSds871cQwVLa_uwzAHgNe1fbrq0UjEPcMffXd-3TCbnWCw4sAoRfyFgobq4lgOqchx_aFP7wrhm_EPZ1nHJOej5PyPz-cPYOx mbiRLQq2jKWGz-Tq4Ysu3UvjC4AG_7PgRUZufFkJi1q1zfZOZ1BxLHTI3nBlv94QLd7408IA6eXZs2B_4NR8SdJ2Mz2FXwQ7a94Ew&state=YJoRG09FnI p8jfA1fv1gRwKpV8dXsL21M74yfsFB69M (https://sites.google.com/view/cs110-final-project/home?code=AQAKDB1BLnDXqY05zh8WaB7 W1ZhqXzm gTJsc85RfqoCSbvobdkPUBLIibSX2S1Th0VVgZY00yxUjEJRMXPel_pQegDpG4KQHxJUcDUB_DsbfnECL8qOVwis31FG4dxXbZyi5u0SnzPD42 NF0cakywxSibOgGaLdRT0d0LNLZQ3vq4n2ajLr-0YS2ZkC603rminveMx6dePloLI35494cdY7z6MPVJd5KU01_VeQey_QasEKoxTBPrU6FxOfpPCULMjx 01E77LW5FiV9VS3-WTmk89hcNINF8awVA12n8Gu8ynhBE3IUQUloXFEEIhnvufW27IH4y73twX3uo2DHW-eY7xPonq-WnDX76pZOV0iKQ_XhF_23ALTA7h 0bEPJLig2o_pds0BTSWFpyNQ-z9RuX-jq7y0vR_C6X0jUlF5mkB41nHbGq_N7zwDEEoNXCZ6CrjdLUdBgpMzcAZkuGuHanj95yYa2Tna6QFIa7VV2u7tic 8kELrLCp1A-NROD_PKviC4QLza81HSds871cQwVLa_uwzAHgNe1fbrq0UjEPcMffXd-3TCbnWCw4sAoRfyFgobq4lgOqchx_aFP7wrhm_EPZ1nHJOej5Py Pz-cPYOxmbiRLQq2jKWGz-Tq4Ysu3UvjC4AG_7PgRUZufFkJi1q1zfZOZ1BxLHTI3nBlv94QLd7408IA6eXZs2B_4NR8SdJ2Mz2FXwQ7a94Ew&state=YJ oRG09FnIp8jfA1fv1gRwKpV8dXsL21M74yfsFB69M)
Searching for "Passionfruit" by Drake...
Done.
Searching for "Relaxa!" by Haikaiss...
Done.
Searching for "Weird Fishes/ Arpeggi" by Radiohead...
Done.
Searching for "Fake Plastic Trees" by Radiohead...
Done.
Searching for "Pull Me Under" by Dream Theater...
Done.
Searching for "Whole Lotta Love" by Led Zeppelin...
Done.
Searching for "Linda" by Projota...
Done.
Searching for "Passarinhos" by Emicida...
Done.
Searching for "Here I Go Again" by Whitesnake...
Done.
Searching for "Fly Me To The Moon" by Frank Sinatra...
Done.
[(*'my'*, 36), (*'no'*, 32), (*'love'*, 24), (*'não'*, 21), (*'again'*, 15), (*'whole'*, 15), (*'here'*, 15), (*'você'*, 14), (*'all'*, 14), (*'go'*, 14), (*'want'*, 13), (*'out'*, 12), (*'meu'*, 12), (*'assim'*, 12), (*'wears'*, 12), (*'pull'*, 12), (*'under'*, 12), (*'baby'*, 12), (*'lotta'*, 12), (*'up'*, 11)]

In [21]:

1

songTable.repetitions('need')

Out[21]: 5

LOs and HCs

LOs

- #DataStructes - Used, modified, and explained different data-structures, providing technical information about them. Also compared 2 different types of data-structures for the same application, using both an experimental and analytical approach.
- #RandomizationTechniques - Used randomization techniques when doing the test to compare the data structures in order to get a less biased bank of words and reduce the chance of skewing my results in favor of one data-structure over the other.
- #ComputationalSolutions - Effectively broke down the different steps of the problem, formulating different algorithmic solutions, in order to later combine them and design an elegant application.
- #ComputationalCritique - Identified strengths and weaknesses of the different data-structures and decided which one was the best to use given my scenario.
- #ComplexityAnalysis - Carefully analyzed the asymptotic behavior of my solutions both analytically and experimentally, using this result to decide which one was more effective.
- #PythonPrograming - Effectively used libraries and APIs, combined with my own solutions and codes to design a python application to my problem.

HCs

- #communicationdesign - Designed a code that is easy to use and provided a clear explanation on how to use it and what to do if it does not work, using memes and informal language as a way for the user to better engage with the application.

#selfawareness - On my initial proposal I had another optional item that I was planning to do as well. However, I realized that it would be better for me to focus on doing one great and 2 bad. Trying to do the other application on time would have resulted in 2 bad applications instead of 1 good one.

#organization - Clearly organized my code dividing it into sub-parts and providing good comments and explanations that will help anyone that reads it to understand it better.

In []:

▶

1