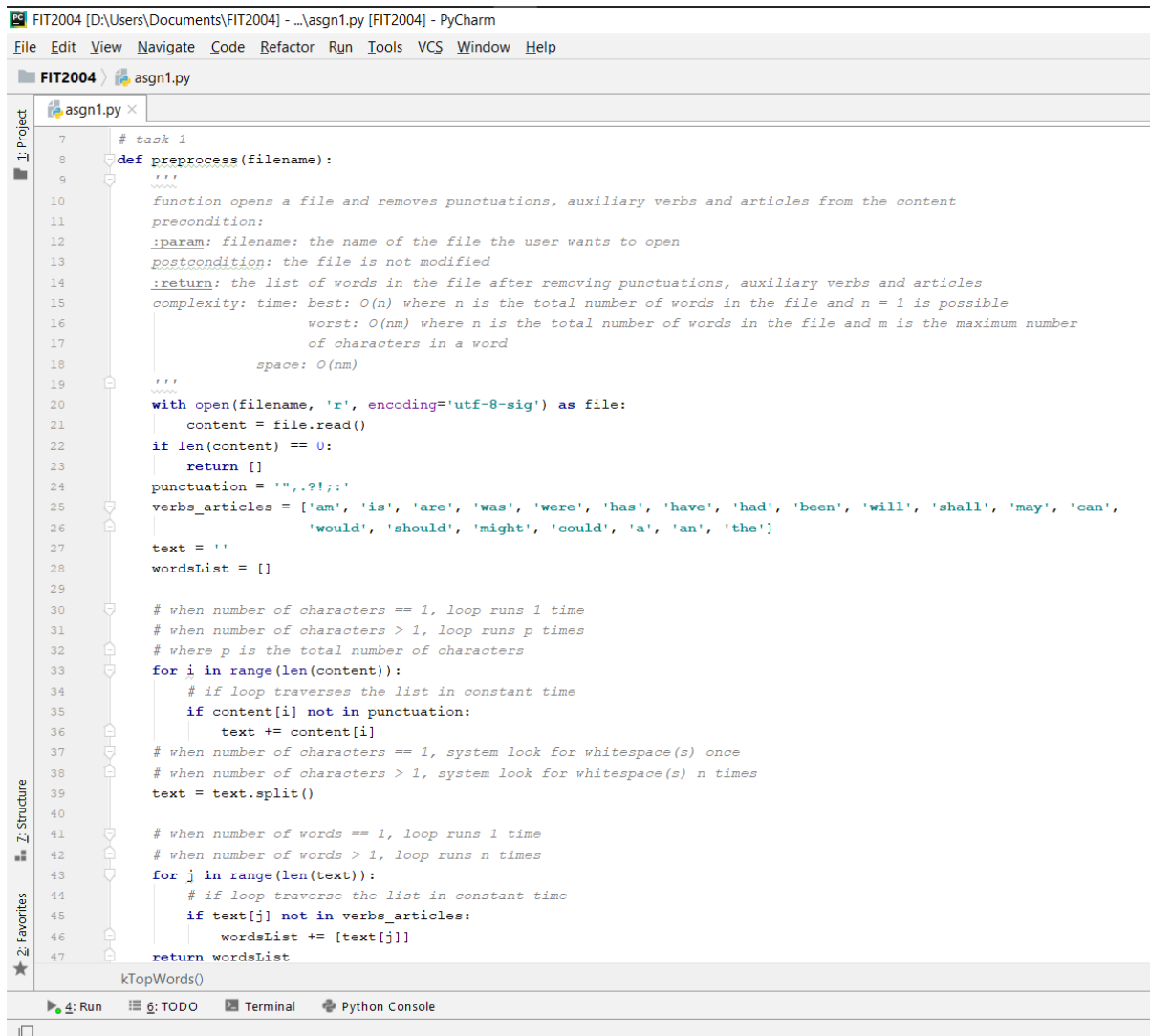


ID: 28390121

Name: Priscilla Tham Ai Ching

FIT2004 Assignment 1: Analysis

TASK 1



```
7 # task 1
8 def preprocess(filename):
9     """
10     function opens a file and removes punctuations, auxiliary verbs and articles from the content
11     precondition:
12     :param: filename: the name of the file the user wants to open
13     postcondition: the file is not modified
14     :return: the list of words in the file after removing punctuations, auxiliary verbs and articles
15     complexity: time:  $O(n)$  where  $n$  is the total number of words in the file and  $n = 1$  is possible
16                 worst:  $O(nm)$  where  $n$  is the total number of words in the file and  $m$  is the maximum number
17                   of characters in a word
18                 space:  $O(nm)$ 
19     """
20     with open(filename, 'r', encoding='utf-8-sig') as file:
21         content = file.read()
22     if len(content) == 0:
23         return []
24     punctuation = ".,?!;:"
25     verbs_articles = ['am', 'is', 'are', 'was', 'were', 'has', 'have', 'had', 'been', 'will', 'shall', 'may', 'can',
26                       'would', 'should', 'might', 'could', 'a', 'an', 'the']
27     text = ''
28     wordsList = []
29
30     # when number of characters == 1, loop runs 1 time
31     # when number of characters > 1, loop runs p times
32     # where p is the total number of characters
33     for i in range(len(content)):
34         # if loop traverses the list in constant time
35         if content[i] not in punctuation:
36             text += content[i]
37     # when number of characters == 1, system look for whitespace(s) once
38     # when number of characters > 1, system look for whitespace(s) n times
39     text = text.split()
40
41     # when number of words == 1, loop runs 1 time
42     # when number of words > 1, loop runs n times
43     for j in range(len(text)):
44         # if loop traverse the list in constant time
45         if text[j] not in verbs_articles:
46             wordsList += [text[j]]
47     return wordsList
```

kTopWords()

Run TODO Terminal Python Console

Task 1 requires a function *preprocess* which will take a filename of the input file and return a list of words in the file after removing auxiliary verbs, punctuation and articles. However, the file itself should not be modified.

The in-built function *open()* reads a file as it takes a filename with its extension which we can assume it as an operation that cost constant time complexity represented in big O notation as $O(1)$.

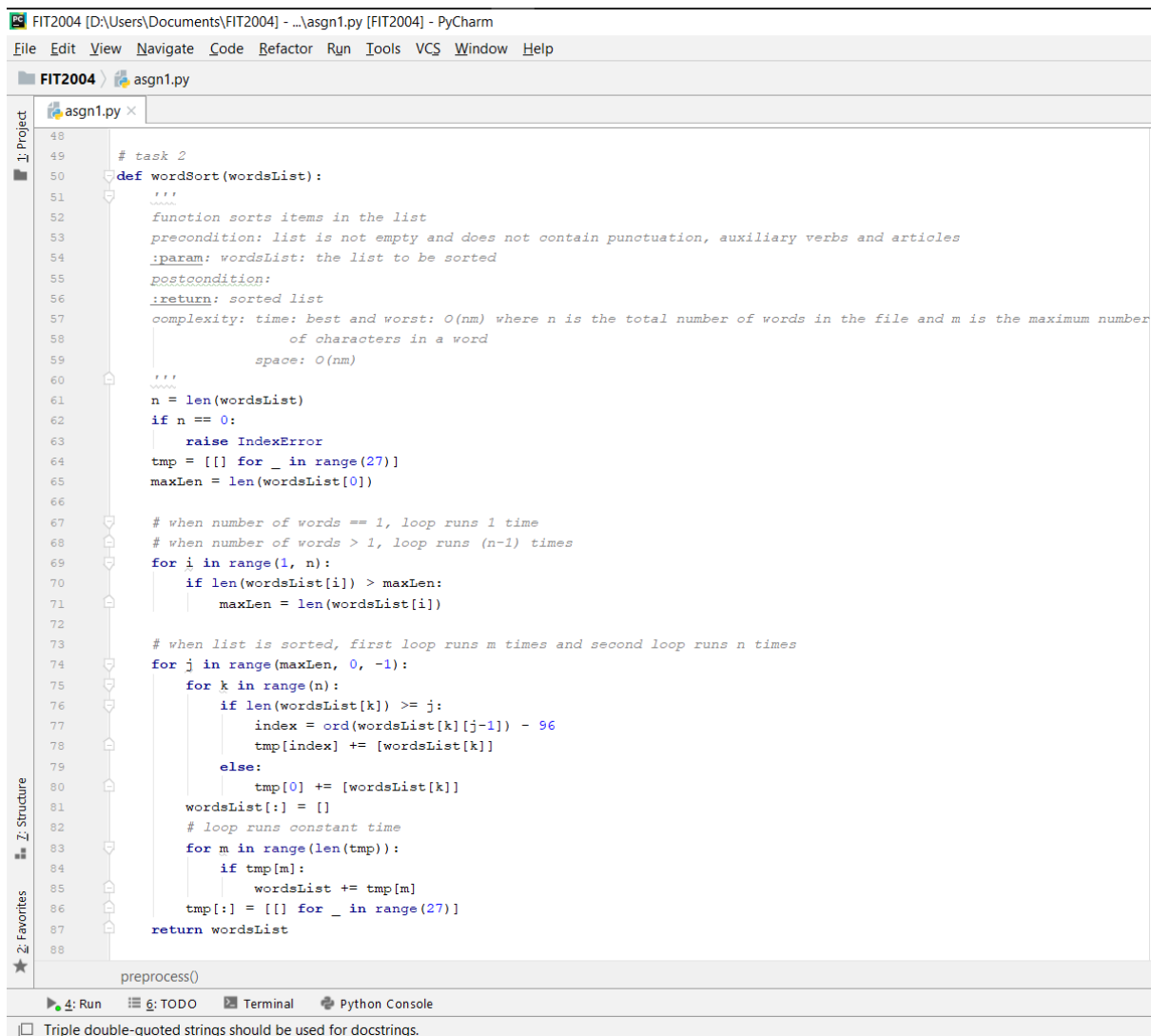
The first *for* loops through the characters in the file denoted as p . Then, the *if-statement* is a 'loop' to compare strings in the file with items in the list containing punctuation. This comparison is constant as both the string and the item are always of length = 1.

This concatenated characters are then separated by whitespace(s) using the in-built function `split()`. The function looks for whitespace(s) between the words (end of each word). So, consider n to be the number of words then, the function looks for it n times. The total time complexity so far is $p + n + 7$.

The second *for* loops through the words in the file, again, denoted as n . Then, the *if-statement* is a 'loop' to compare strings in the file with items in the list containing auxiliary verbs and articles. So, consider m to be the maximum number of characters in a word then, the comparison cost is m where m equals the maximum 'loop'. Therefore, the total time complexity is $(nm) + n + p + 8$, ignoring the constant we get, $(nm) + n + p$. We know that $nm > p > n$, hence, $O(nm)$ is the worst time complexity.

Ignoring the space taken by the list of punctuation and the list of auxiliary verbs and articles for it is constant, the list of characters and the list of words are affected by the file input. Appending to the list of characters cost p and appending to the list of words cost n so total space complexity is $p + n$. However, since strings are list itself, it cost m when appending to the list of words where m (maximum number of characters in a word) equals to the maximum size of list inside the list of words. Hence, $O(nm)$ is the worst space complexity since $nm > p > n$.

TASK 2



```
48
49 # task 2
50 def wordSort(wordsList):
51     """
52     function sorts items in the list
53     precondition: list is not empty and does not contain punctuation, auxiliary verbs and articles
54     :param wordsList: the list to be sorted
55     :postcondition:
56     :return: sorted list
57     complexity: time: best and worst: O(nm) where n is the total number of words in the file and m is the maximum number
58                 of characters in a word
59                 space: O(nm)
60     """
61     n = len(wordsList)
62     if n == 0:
63         raise IndexError
64     tmp = [[] for _ in range(27)]
65     maxLen = len(wordsList[0])
66
67     # when number of words == 1, loop runs 1 time
68     # when number of words > 1, loop runs (n-1) times
69     for i in range(1, n):
70         if len(wordsList[i]) > maxLen:
71             maxLen = len(wordsList[i])
72
73     # when list is sorted, first loop runs m times and second loop runs n times
74     for j in range(maxLen, 0, -1):
75         for k in range(n):
76             if len(wordsList[k]) >= j:
77                 index = ord(wordsList[k][j-1]) - 96
78                 tmp[index] += [wordsList[k]]
79             else:
80                 tmp[0] += [wordsList[k]]
81         wordsList[:] = []
82         # loop runs constant time
83         for m in range(len(tmp)):
84             if tmp[m]:
85                 wordsList += tmp[m]
86         tmp[:] = [[] for _ in range(27)]
87     return wordsList
88
preprocess()
```

Task 2 requires a function *wordSort* which will take a list of words and return an alphabetically sorted list of words. The algorithm of radix sort is universal so, we would skip on further elaboration but the implementation to unequal length of string here, is done by comparing the length rather than adding space or special characters behind to make them all equal length. For strings that we cannot compare alphabetically (length of the string is < the length of the longest string), those strings are placed in the first position of a temporary list to be compared later when the point of comparison is not the last alphabet of the longest string.

The first *for* loops through the words in the list denoted as n . Then, the *if-statement* compares the length of one word to another with the assumption that the first word in the list is the longest. This uses the in-built function *len()* which we can assume it as an operation that cost constant time complexity represented in big O notation as $O(1)$. The total time complexity so far is $(n-1) + 5$.

The second *for* loops length-of-longest-word times. So, consider m to be the maximum number of characters in a word then, the loop runs m times. There are two *for* loops within. The former loops through the words in the list, again, denoted as n . Then, the *if-statement* compares the length of a word to the counter of the outer loop which starts from the length of

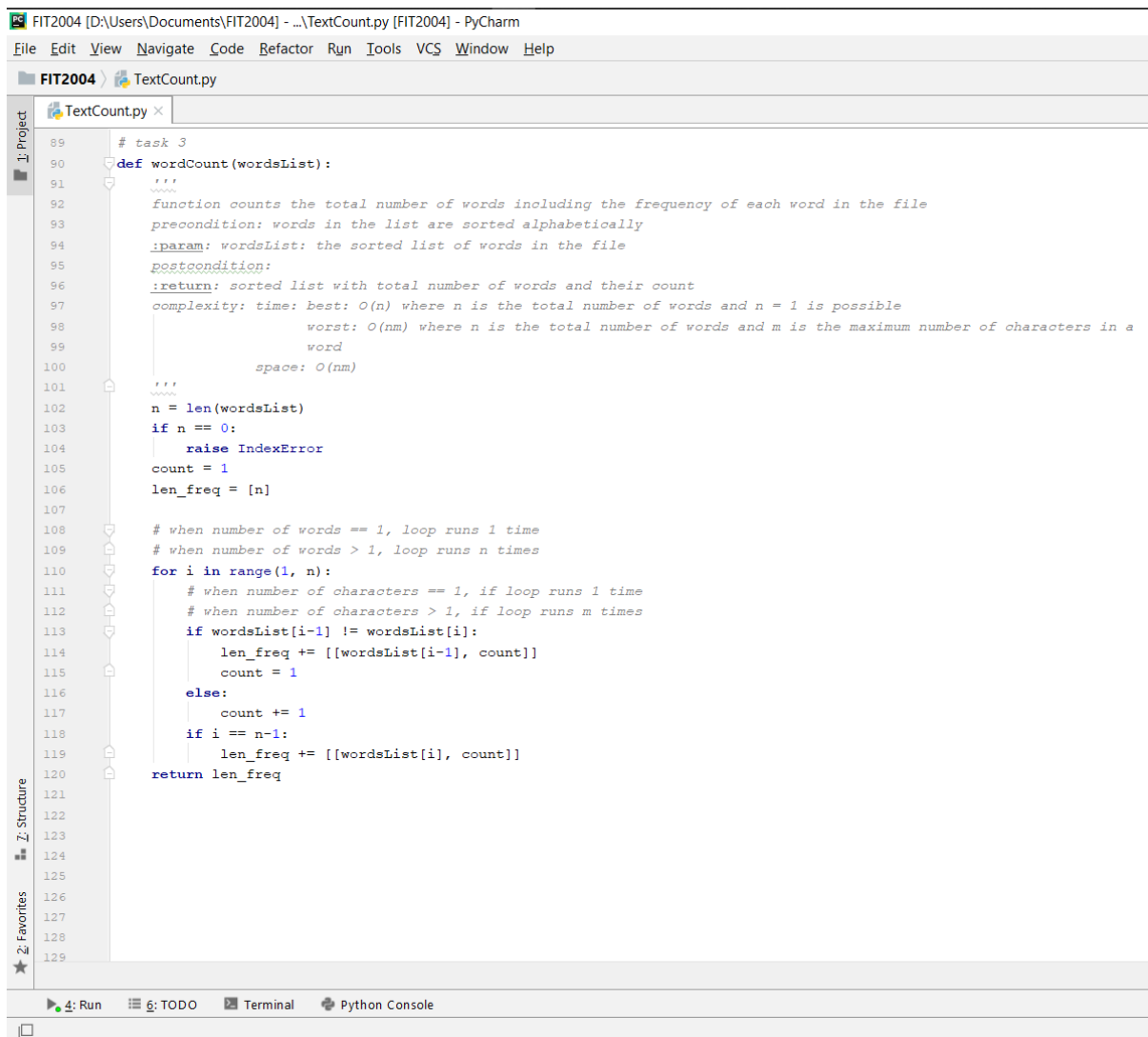
the longest word. This uses the in-built function *len()* which we can assume it as an operation that cost constant time complexity represented in big O notation as $O(1)$. The condition within calculates the index of the alphabet ($a = 1$ and so on) to fill it into the temporary list. This uses the in-built function *ord()* which we can assume it as an operation that cost constant time complexity represented in big O notation as $O(1)$.

Once the former loop is over, the list of words will then be sliced to remove the items in the list. This operation is in constant time as a complete slice is performed where list traversal is not done. The latter loops through the temporary list of size 27. This is also an operation of constant time as the size is constant referring to the English alphabets. Once the latter loop is over, the temporary list is also sliced to remove the items in the list and the time complexity is also as explained for previous slicing. The total time complexity of both loops inside is $n + 27 + 7 = n + 34$.

Therefore, the total time complexity is $(n-1) + 5 + m(n + 34) = nm + n + 34m + 4$, ignoring the constant we get, $nm + n + 34m$. We know that $nm > m > n$, hence, $O(nm)$ is the worst time complexity.

Ignoring the space taken by the temporary list for it is constant, the list of words is affected by the file input. Since list is mutable, the list of words passed as an argument here is a reference to the returned list from task 1. Hence, as explained from task 1, the worst space complexity is $O(nm)$.

TASK 3



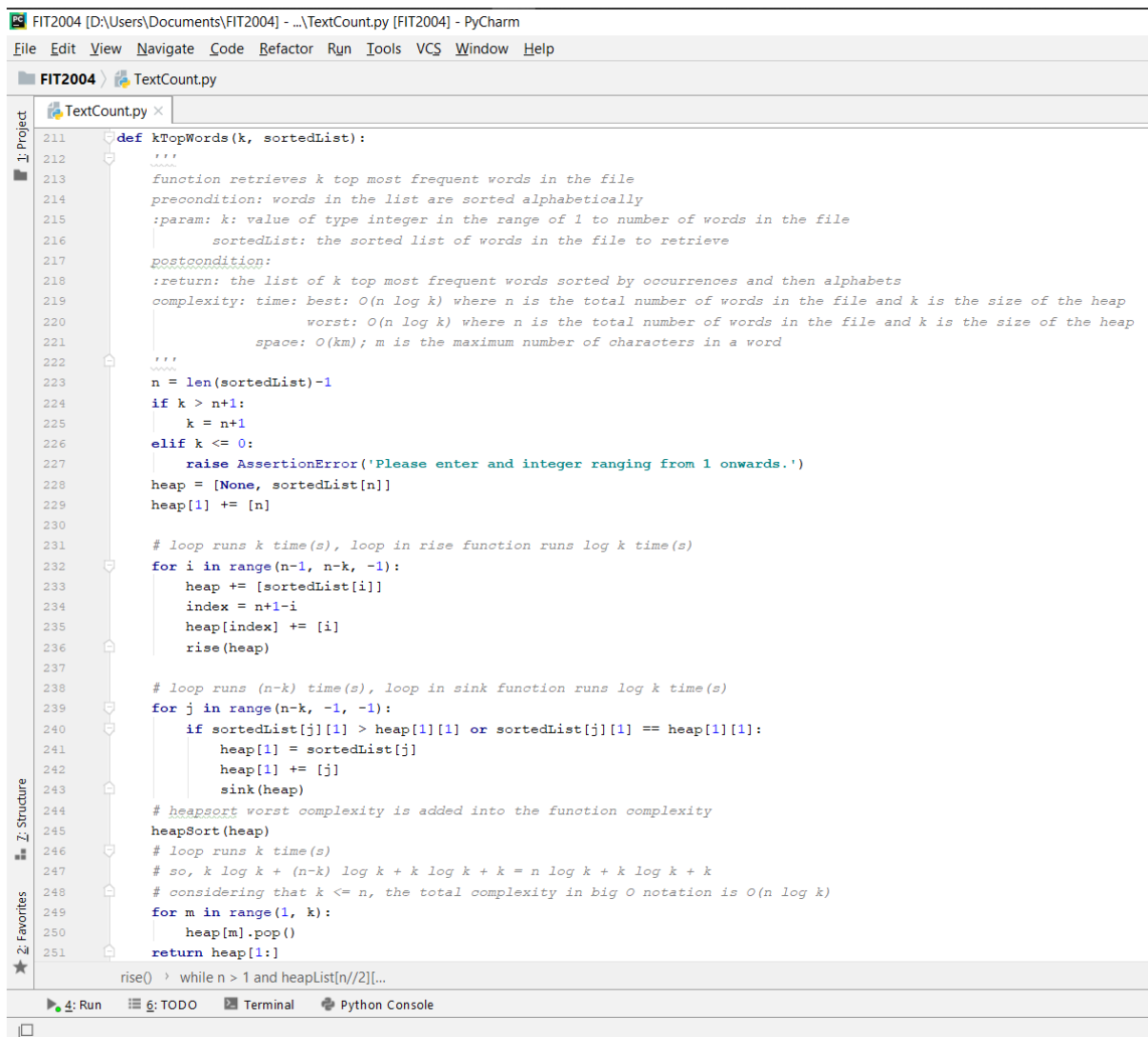
```
89 # task 3
90 def wordCount(wordsList):
91     """
92     function counts the total number of words including the frequency of each word in the file
93     precondition: words in the list are sorted alphabetically
94     :param wordsList: the sorted list of words in the file
95     postcondition:
96     :return: sorted list with total number of words and their count
97     complexity: time: best: O(n) where n is the total number of words and n = 1 is possible
98                worst: O(nm) where n is the total number of words and m is the maximum number of characters in a
99                word
100                space: O(nm)
101     """
102     n = len(wordsList)
103     if n == 0:
104         raise IndexError
105     count = 1
106     len_freq = [n]
107
108     # when number of words == 1, loop runs 1 time
109     # when number of words > 1, loop runs n times
110     for i in range(1, n):
111         # when number of characters == 1, if loop runs 1 time
112         # when number of characters > 1, if loop runs m times
113         if wordsList[i-1] != wordsList[i]:
114             len_freq += [[wordsList[i-1], count]]
115             count = 1
116         else:
117             count += 1
118         if i == n-1:
119             len_freq += [[wordsList[i], count]]
120     return len_freq
```

Task 3 requires a function *wordCount* which will take an alphabetically sorted list of words and return a list with two values: (a) the total number of words and (b) a list of words with their count. However, the list of words should remain sorted.

The *for* loops through the words in the list denoted as n . Then, the *if-statement* is a ‘loop’ to compare previous word to the next word in the list. So, consider m to be the maximum number of characters in a word then, the comparison cost is m where m equals the maximum ‘loop’. Therefore, total time complexity is $nm + 6$, ignoring the constant we get, nm , hence, the worst time complexity is $O(nm)$.

The list created to return is affected by the file input and is of size $(n+1)$. Appending to the list will cost n so total space complexity is $n+1$ (one space for the total number of words). However, since strings are list itself, it cost m when appending to the list where m (maximum number of characters in a word) equals to the maximum size of list inside the list of words. Hence, $O(nm)$ is the worst space complexity.

TASK 4



```
211 def kTopWords(k, sortedList):
212     """
213     function retrieves k top most frequent words in the file
214     precondition: words in the list are sorted alphabetically
215     :param k: value of type integer in the range of 1 to number of words in the file
216     |         sortedList: the sorted list of words in the file to retrieve
217     postcondition:
218     :return: the list of k top most frequent words sorted by occurrences and then alphabets
219     complexity: time: best:  $O(n \log k)$  where n is the total number of words in the file and k is the size of the heap
220     |         worst:  $O(n \log k)$  where n is the total number of words in the file and k is the size of the heap
221     |         space:  $O(km)$ ; m is the maximum number of characters in a word
222     """
223     n = len(sortedList)-1
224     if k > n+1:
225         k = n+1
226     elif k <= 0:
227         raise AssertionError('Please enter and integer ranging from 1 onwards.')
228     heap = [None, sortedList[n]]
229     heap[1] += [n]
230
231     # loop runs k time(s), loop in rise function runs log k time(s)
232     for i in range(n-1, n-k, -1):
233         heap += [sortedList[i]]
234         index = n+1-i
235         heap[index] += [i]
236         rise(heap)
237
238     # loop runs (n-k) time(s), loop in sink function runs log k time(s)
239     for j in range(n-k, -1, -1):
240         if sortedList[j][1] > heap[1][1] or sortedList[j][1] == heap[1][1]:
241             heap[1] = sortedList[j]
242             heap[1] += [j]
243             sink(heap)
244     # heapSort worst complexity is added into the function complexity
245     heapSort(heap)
246     # loop runs k time(s)
247     # so,  $k \log k + (n-k) \log k + k \log k + k = n \log k + k \log k + k$ 
248     # considering that  $k \leq n$ , the total complexity in big O notation is  $O(n \log k)$ 
249     for m in range(1, k):
250         heap[m].pop()
251     return heap[1:]
```

Task 4 requires a function *kTopWords* which will take a value *k* and an alphabetically sorted list of words with their frequencies and return a list that contains the *k* number of top-most frequent words. However, if the count of two words are the same, the word comes earlier in the sorted list will have higher priority over another. The algorithm to create a min heap and sort the heap is universal so, we would skip further elaboration but the implementation to stabilise the heap before and after sorting here, is done with the concept of priority queue. Position of the words in the sorted list is appended to it whenever the word at the root of the heap is replaced. This position value is what determines the priority of words to let remain/be replaced in the heap. Also, creating the min heap is done from the end of the list and sorting the heap is done from the bottom of the heap. This is to simplify replacement of words of higher priority by words of higher occurrences and the replacement of words of lower priority by those of higher priority when their occurrences are the same.

```
195 def rise(heapList):
196     """
197     function moves item up the heap accordingly
198     precondition:
199     :param heapList: the heap array
200     postcondition: heap array is stable
201     :return:
202     complexity: time: best: O(log k) where k is the size of the heap
203                worst: O(log k) where k is the size of the heap
204     """
205     n = len(heapList)-1
206     # loop runs log k time(s)
207     while n > 1 and heapList[n//2][1] > heapList[n][1]:
208         heapList[n//2], heapList[n] = heapList[n], heapList[n//2]
209         n //= 2
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
```

kTopWords()

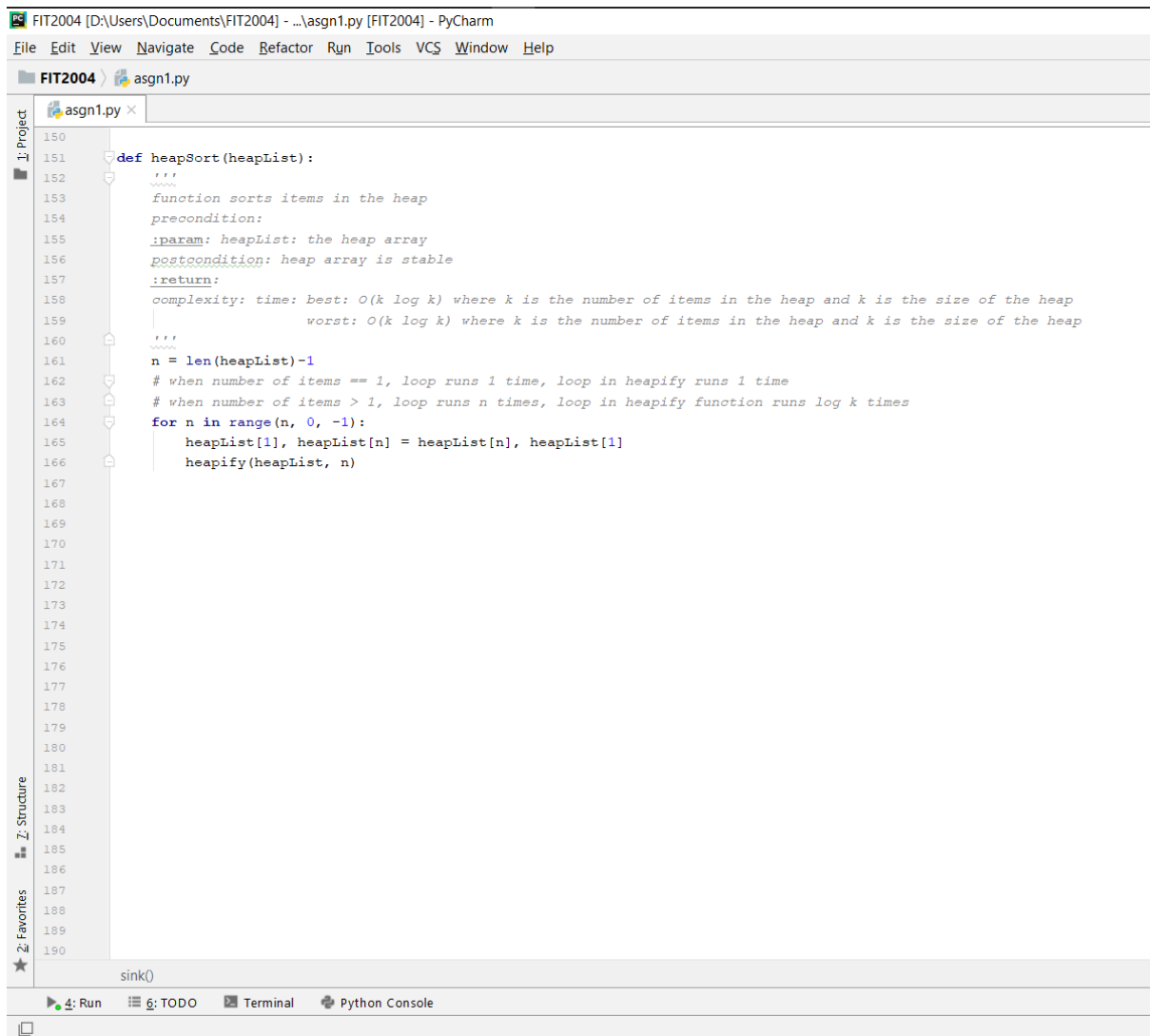
Run TODO Terminal Python Console

The first *for* loops through the $k-1$ words in the sorted list to append the $k-1$ items into the heap. The function *rise* is then called to ensure the heap is a min heap. In the function *rise*, the *while* loop compares item in the heap at a certain position(left/right child) to another item at position $n//2$ (parent node)(integer division). The time complexity $k \log k$ so far is the total after both the loop and the *rise* function.

```
168 def sink(heapList):
169     """
170     function moves item down the heap accordingly
171     precondition:
172     :param: heapList: the heap array
173     postcondition: heap array is stable
174     :return:
175     complexity: time: best: O(log k) where k is the size of the heap
176                worst: O(log k) where k is the size of the heap
177     """
178     root = 1
179     # loop runs log k time(s)
180     while 2*root <= len(heapList)-1:
181         child = 2*root
182
183         if child+1 < len(heapList) and (heapList[child+1][1] < heapList[child][1] or
184             (heapList[child][1] == heapList[child+1][1] and
185             heapList[child][2] < heapList[child+1][2])):
186             child += 1
187
188         if heapList[root][1] > heapList[child][1] or \
189             (heapList[child][1] == heapList[root][1] and heapList[child][2] > heapList[root][2]):
190             heapList[root], heapList[child] = heapList[child], heapList[root]
191         else:
192             return
193         root = child
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
```

The second *for* loops through the remaining words in the sorted list. So, consider n to be the total number of words then, the loop runs $n-k$ times. Then, the *if-statement* compares the occurrence of the word to the occurrence of the word at the root of the heap. This comparison is constant as both is integer. The function *sink* is then called to ensure the heap is a min heap. In the function *sink*, the *while* loop compares item in the heap at a certain position(parent node) to another item at position $*2$ and/or position $*2+1$ (left and right child respectively). The time complexity $k \log k + (n-k) \log k$ so far is the total after both previous loop and this loop.

Adding to the implementation, there are extra conditions added to the function compared to a normal *sink* function. The choice of child when both left and right occurrences of words is equal depends on the position of the words in the sorted list. The choice of moving the root item down the heap by swapping the child is also determined the same way. This is to ensure the lower prioritised words are first replaced.



```
150
151 def heapSort(heapList):
152     """
153     function sorts items in the heap
154     precondition:
155     :param: heapList: the heap array
156     postcondition: heap array is stable
157     :return:
158     complexity: time: best:  $O(k \log k)$  where  $k$  is the number of items in the heap and  $k$  is the size of the heap
159                worst:  $O(k \log k)$  where  $k$  is the number of items in the heap and  $k$  is the size of the heap
160     """
161     n = len(heapList)-1
162     # when number of items == 1, loop runs 1 time, loop in heapify runs 1 time
163     # when number of items > 1, loop runs n times, loop in heapify function runs log k times
164     for n in range(n, 0, -1):
165         heapList[1], heapList[n] = heapList[n], heapList[1]
166         heapify(heapList, n)
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
sink()
```

That said, the heap now has all the necessary words and their frequency. It is pass as reference to the function *heapSort*. In the function *heapSort*, the *for* loops through the words in the heap, again, denoted as *k* and pass it as reference to the function *heapify* to swap items around until it is sorted. The *k* mentioned is the same as the *k* in previous function(s).

```
120 # task 4
121 def heapify(heapList, index):
122     """
123     function sorts items in the heap by sub-heap tree
124     precondition:
125     :param heapList: the heap array
126     postcondition: heap array is stable
127     :return:
128     complexity: time: best: O(log k) where k is the size of the heap
129                worst: O(log k) where k is the size of the heap
130     """
131     root = 1
132     # loop runs log k time(s)
133     while 2*root <= index:
134         smallest = root
135         left = 2*root
136         right = 2*root+1
137
138         if left < index and (heapList[left][1] < heapList[smallest][1] or
139                             (heapList[left][1] == heapList[smallest][1] and heapList[left][2] > heapList[smallest][2])):
140             smallest = left
141
142         if right < index and (heapList[right][1] < heapList[smallest][1] or
143                             (heapList[right][1] == heapList[smallest][1] and heapList[right][2] > heapList[smallest][2])):
144             smallest = right
145
146         if smallest == root:
147             return
148         heapList[smallest], heapList[root] = heapList[root], heapList[smallest]
149         root = smallest
150
151
152
153
154
155
156
157
158
159
160
```

In the function *heapify*, the *while* loop compares item in the heap at a certain position(parent node) to another item at position*2 and/or position*2+1(left and right child respectively). The time complexity $k \log k + (n-k) \log k + k \log k = n \log k + k \log k$ so far is the total.

Adding to the implementation, there are extra conditions added to the function compared to a normal *heapify* function. The choice of child when both left and right occurrences of words is equal depends on the position of the words in the sorted list. The choice of swapping the root item with the child item is also determined the same way. This is to ensure the higher prioritised words are placed in the front of the lower prioritised words if their occurrences are equal.

Once the heap is sorted, the third for loops through the items in the heap, *k*, to remove unnecessary items/details before returning. This uses the in-built function *pop()* which we can assume it as an operation that cost constant time complexity represented in big O notation as $O(1)$ as the item is in the last position; no list traversal is needed. Therefore, the total time complexity is $n \log k + k \log k + k = (n+k) \log k + k$. We know that $k \leq n$ so, when $k = n$, $2n \log k > n$, hence, the worst time complexity is $O(n \log k)$.

The list (heap) created to return is affected by the *k* value of the user input. Appending to the list will cost *k* so total space complexity is *k*. However, since strings are list itself, it cost *m*

when appending to the list where m (maximum number of characters in a word) equals to the maximum size of list inside the list/heap. Hence, $O(nm)$ is the worst space complexity.