

Contents

1	Basic Test Results	2
2	README	3
3	PathScanner.py	5
4	WordExtractor.py	8
5	WordTracker.py	10

1 Basic Test Results

```
1 Starting tests...
2 Thu Jan 15 19:58:36 IST 2015
3 4056f36443af52b9040ef228e27db7b5c4009071 -
4
5
6 -rw-r--r-- ransha/stud    2380 2015-01-15 19:55 README
7 -rwxr--r-- ransha/stud    5091 2015-01-15 19:50 WordTracker.py
8 -rwxr--r-- ransha/stud    2774 2015-01-15 18:24 WordExtractor.py
9 -rwxr--r-- ransha/stud    5692 2015-01-15 18:15 PathScanner.py
10
11 Testing README...
12 Done testing README...
13
14 Running presubmit tests...
15 result_code    wtrck_example    4    1
16 result_code    ps_example    5    1
17 result_code    wext_example    3    1
18 Done running presubmit tests
19
20 Tests completed
```

2 README

```
1  ransha
2  203781000
3  Ran Shaham
4
5  =====
6  = README for ex10: Iterators and files =
7  =====
8
9  =====
10 = Description: =
11 =====
12
13 In this exercise we implemented 3 programs: one for reading all words in a
14 certain file, one that determines if a given word is in a certain list and the
15 last is for searching files containing a given words list in the file system.
16 All of those were implemented using iterator, classes, files i/o (actually
17 only i...) and recursion.
18
19 =====
20 = List of submitted files: =
21 =====
22
23 README          This file
24
25 WordExtractor.py  An iterator that extracts every word from a file
26                  that contains text.
27 WordTracker.py   A class that takes a 'dictionary' (word list) as and
28                  contains methods that sort and search the list.
29 PathScanner.py   Contains methods that use the other files (from this
30                  exercise) and search for file in the file system.
31
32
33 =====
34 = NOTE =
35 =====
36
37 - Part 2 - WordTracker.py -
38
39 __contains__() method:
40 In order to use binary search in the __contains__ function the list of words
41 needs to be sorted. Therefore, I implemented 'Merge Sort' in order to
42 efficiently sort the list (_sort_dict() function that runs on __init__).
43 The instructions specified the implementation need to be efficient in terms
44 of runtime complexity, so the merge sort, in my opinion, is a good choice.
45 Although it uses recursion and a relatively large amount of memory, it runs
46 in  $O(n \log n)$  which is better (or at least as good) than other sorts.
47
48 encounter() method:
49 this function uses the __contains__ method that runs in  $O(n)$  and returns its
50 value. In the __contains__ method the found words list is updated before it
51 returns True, so that the encountered_all() function will run fast.
52 Therefore, worst-case runtime complexity will still be  $O(n)$ .
53
54 NOTE - I think binary search has runtime complexity of  $O(\log n)$  not like
55 mentioned in the instructions ( $O(n)$ ).
56
57
58 encountered_all() method:
59 Goes over every item in the found words list (max n items) to check if one of
```

```
60 the words wasn't found. therefore: worst-case runtime complexity of  $O(n)$ .
61
62 reset() method:
63 Sets a new list of 'False' values, thus resetting the search for words.
64 Runtime complexity of  $O(n)$ .
```

3 PathScanner.py

```
1 import os
2 from WordExtractor import *
3 from WordTracker import *
4
5 class PathIterator:
6     """
7     An iterator which iterates over all the directories and files
8     in a given path (note - in the path only, not in the
9     full depth). There is no importance to the order of iteration.
10    """
11    def __init__(self, path):
12        self.items = os.listdir(path)
13        self.index = 0
14        self.path = path
15        self.SEPERATOR = '/'
16
17    def __iter__(self):
18        return self
19
20    def __next__(self):
21        """
22        Goes over the items list, as given by the os.listdir() function one
23        item at a time.
24        returns the path of a file/dir in the directory as string
25        """
26        if self.index >= len(self.items): raise StopIteration()
27        item = self.items[self.index]
28        self.index += 1
29        return self.path + self.SEPERATOR + item
30
31
32
33 def path_iterator(path):
34     """
35     Returns an iterator to the current path's files and directories.
36     Note - the iterator class is not outlined in the original
37     version of this file - but rather it should be designed
38     and implemented by you.
39     :param path: A (relative or an absolute) path to a directory.
40     It can be assumed that the path is valid and that indeed it
41     leads to a directory (and not to a file).
42     :return: An iterator which returns all the files and directories
43     in the *current* path (but not in the *full depth* of the path).
44     """
45     return PathIterator(path)
46
47 def print_tree(path, sep=' '):
48     """
49     Print the full hierarchical tree of the given path.
50     Recursively print the full depth of the given path such that
51     only the files and directory names should be printed (and not
52     their full path), each in its own line preceded by a number
53     of separators (indicated by the sep parameter) that correlates
54     to the hierarchical depth relative to the given path parameter.
55     :param path: A (relative or an absolute) path to a directory.
56     It can be assumed that the path is valid and that indeed it
57     leads to a directory (and not to a file).
58     :param sep: A string separator which indicates the depth of
59     current hierarchy.
```

```

60     """
61
62     def recursive_tree(path, sep, depth):
63         """
64         The recursive function that goes over the path.
65         Gets a path of a DIRECTORY, prints every file in it and recursively
66         enters directories in it.
67         :param path: A relative or an absolute path to a directory.
68         :type path: string
69         :param sep: a string which indicates the depth of current hierarchy.
70         :type sep: string
71         :param depth: The current directory depth (distance from input folder)
72         :type depth: int
73         """
74         # an index for the 'absolute' name of file (or dir) inside 'path'
75         path_index = len(path) + 1
76         items = path_iterator(path)
77         for item in items:
78             print(sep * depth + item[path_index:])
79             if os.path.isdir(item):
80                 recursive_tree(item, sep, depth + 1)
81         # starts printing with 0 depth.
82         recursive_tree(path, sep, 0)
83
84
85
86
87
88
89     def file_with_all_words(path, word_list):
90         """
91         Find a file in the full depth of the given path which contains
92         all the words in word_list.
93         Recursively go over the files in the full depth of the given
94         path. For each, check whether it contains all the words in
95         word_list and if so return it.
96         :param path: A (relative or an absolute) path to a directory.
97         In the full path of this directory the search should take place.
98         It can be assumed that the path is valid and that indeed it
99         leads to a directory (and not to a file).
100         :param word_list: A list of words (of strings). The search is for
101         a file which contains this list of words.
102         :return: The path to a single file which contains all the
103         words in word_list if such exists, and None otherwise.
104         If there exists more than one file which contains all the
105         words in word_list in the full depth of the given path, just one
106         of theses should be returned (does not matter which).
107         """
108
109         def recursive_search(path, word_tracker):
110             """
111             Recursive search for the file containing all of the words in
112             the input word_list, using the class WordTracker.
113             :param path: Relative or absolute path of a directory. a string.
114             :param word_tracker: a WordTracker object that holds the word list
115             that we are looking for in the file system.
116             """
117             items = path_iterator(path)
118
119             for item in items:
120                 if os.path.isfile(item):
121
122                     try: # to avoid file system errors
123                         words = WordExtractor(item)
124                         for word in words: # checks every word in the file
125                             word_tracker.encounter(word)
126                         if word_tracker.encountered_all():
127                             return item # if desired file was found stop searching

```

```
128         else:
129             word_tracker.reset()
130     except: # will be reached if the file wasn't meant to be read.
131         # if so, resets the tracker and continues the search.
132         word_tracker.reset()
133
134     elif os.path.isdir(item): # recursive search in directories
135         recursive_search(item, word_tracker)
136
137     return None # reached if no file was found.
138 # let the search begin.
139 word_tracker = WordTracker(word_list)
140 return recursive_search(path, word_tracker)
```

4 WordExtractor.py

```
1  #!/usr/bin/env python3
2
3
4  class WordExtractor(object):
5      """
6      This class should be used to iterate over words contained in files.
7      The class should maintain space complexity of  $O(1)$ ; i.e, regardless
8      of the size of the iterated file, the memory requirements of a class
9      instance should be bounded by some constant.
10     To comply with the space requirement, the implementation may assume
11     that all words and lines in the iterated file are bounded by some
12     constant, so it is allowed to read words or lines from the
13     iterated file (but not all of them at once).
14     """
15
16     def __init__(self, filename):
17         """
18         Initiate a new WordExtractor instance whose *source file* is
19         indicated by filename.
20         :param filename: A string representing the path to the instance's
21         *source file*
22         """
23
24         self.NEW_LINE = '\n'
25         self.file = open(filename)
26         self.line = self.file.readline() # the line string
27         self.splitted = self.line.split() # list of words in the line
28         self.curr_index = 0
29         self.end_of_file = False
30
31
32     def __iter__(self):
33         """
34         Returns an iterator which iterates over the words in the
35         *source file* (i.e - self)
36         :return: An iterator which iterates over the words in the
37         *source file*
38         """
39
40         return self
41
42     def _next_line(self):
43         """
44         Reads next line in the file and saves it to the class.
45         """
46         self.line = self.file.readline()
47         self.splitted = self.line.split()
48         self.curr_index = 0
49
50     def _bad_line(self, line):
51         """
52         Checks if a certain line contains no readable words (string).
53         :param line: a string that represents a line in the file.
54         :return: >True if the line is empty or if the end of line is reached
55                 >False if the line is good or it is the end of file.
56         """
57
58         self.end_of_file = len(line) == 0
59         empty_line = len(self.splitted) == 0
```



```

60     end_of_line = self.curr_index >= len(self.splitted)
61
62     # the end of file line isn't a bad line.
63     return (empty_line or end_of_line) and (not self.end_of_file)
64
65 def __next__(self):
66     """
67     Make a single word iteration over the source file.
68     :return: A word from the file.
69     """
70
71     while self._bad_line(self.line):
72         self._next_line() # only read good lines.
73
74     if not self.end_of_file:
75         word = self.splitted[self.curr_index]
76         self.curr_index += 1
77         return word
78     else:
79         self.file.close()
80         raise StopIteration()
81
82

```

5 WordTracker.py

```
1  #!/usr/bin/env python3
2
3
4  class WordTracker(object):
5      """
6      This class is used to track occurrences of words.
7      The class uses a fixed list of words as its dictionary
8      (note - 'dictionary' in this context is just a name and does
9      not refer to the pythonic concept of dictionaries).
10     The class maintains the occurrences of words in its
11     dictionary as to be able to report if all dictionary's words
12     were encountered.
13     """
14
15     def __init__(self, word_list):
16         """
17         Initiates a new WordTracker instance.
18         :param word_list: The instance's dictionary.
19         """
20         self.dict = word_list[:]          # sets a NEW list
21         self._sort_dict()                 # sorts it
22         self.found = [False] * len(self.dict) # list of booleans that flags
23         # self.found_idx = 0              # found words.
24
25     def _sort_dict(self):
26         """
27         Uses merge sort algorithm and python string comparison to sort the
28         given list of words. Uses recursion.
29         :return: a sorted list of words.
30         """
31     def merge_sort(lst):
32         """
33         the recursive sort.
34         """
35         size = len(lst)
36         if size <= 1:
37             return lst # sorted list (contains 1 or less elements)
38         middle = size // 2
39         left, right = merge_sort(lst[:middle]), merge_sort(lst[middle:])
40         lst = merge(left, right) # the merging
41         return lst # returns a sorted list
42
43     def merge(left, right):
44         """
45         Merges two lists according to its' items' values.
46         Gets two parameters 'left','right' as lists and merges them.
47         :return: merged list.
48         """
49         lst = left + right
50         size = len(lst)
51         l, r = 0, 0
52         for i in range(size):
53             l_val = left[l] if l < len(left) else None
54             r_val = right[r] if r < len(right) else None
55             r_bigger = l_val and r_val and l_val < r_val
56             l_bigger = l_val and r_val and l_val >= r_val
57             if r_bigger or r_val == None:
58                 lst[i] = l_val
59                 l += 1
```

```

60         elif l_bigger or l_val == None:
61             lst[i] = r_val
62             r += 1
63         return lst
64     # calls the recursive sort.
65     self.dict = merge_sort(self.dict)
66
67 def __contains__(self, word):
68     """
69     Check if the input word is contained within dictionary.
70     For a dictionary with n entries, this method guarantees a
71     worst-case running time of O(n) by implementing a
72     binary-search.
73     :param word: The word to be examined if contained in the
74     dictionary.
75     :return: True if word is contained in the dictionary,
76     False otherwise.
77     """
78
79     is_in_dict = False
80     lst = self.dict
81     left = 0
82     right = len(lst)
83
84     while left <= right and right <= len(lst) and left >= 0:
85         middle_i = (right+left) // 2
86         middle = lst[middle_i]
87         if middle > word:
88             if middle_i <= 0: break # to avoid infinite loop
89             right = middle_i - 1
90         elif middle < word:
91             if middle_i >= len(lst) - 1: break # same as ^
92             left = middle_i + 1
93         else:
94             is_in_dict = True
95             self.found[middle_i] = is_in_dict
96             break
97         # if the word isn't in the dictionary the middle_i index can
98         # reach the end of the list, and that can cause infinite search.
99     return is_in_dict
100
101
102 def encounter(self, word):
103     """
104     A "report" that the given word was encountered.
105     The implementation changes the internal state of the object as
106     to "remember" this encounter.
107     :param word: The encountered word.
108     :return: True if the given word is contained in the dictionary,
109     False otherwise.
110     """
111     return word in self # uses the __contains__ method
112
113 def encountered_all(self):
114     """
115     Checks whether all the words in the dictionary were
116     already "encountered".
117     :return: True if for each word in the dictionary,
118     the encounter function was called with this word;
119     False otherwise.
120     """
121     for found in self.found:
122         if not found: return found
123     return True
124
125 def reset(self):
126     """
127     Changes the internal representation of the instance such

```

```
128         that it "forget" all past encounters. One implication of
129 such forgetfulness is that for encountered_all function
130 to return True, all the dictionaries' entries should be
131 called with the encounter function (regardless of whether
132 they were previously encountered or not).
133         """
134         self.found = [False] * len(self.dict)
```