# Anagrams

## The problem

Given a words.txt file containing a newline-delimited list of dictionary words, please implement the Anagrams class so that the get_anagrams() method returns all anagrams from words.txt for a given word.

**Bonus requirements:**

- • Optimise the code for fast retrieval

- • Write more tests

- • Thread safe implementation

## General approach

*"An anagram is direct word switch or word play, the result of rearranging the letters of a word or phrase to produce a new word or phrase, using all the original letters exactly once"* ( source: wikipedia )

That means that in order to get all the anagrams for a needed word, we don't need to compare the words theirselves but their ordered representation.

Given two words, word1 and word2

If the ordered characters of word1 are the same that the ordered characters of word2,

Then

word1 and word2 are anagrams.

## Assumtions

- • One given word is anagram of itself.

- • Anagrams are **not** case sensitive so "Star" is an anagram of "Tras".

- • Special caracters as " ' " are considered as regular caracters too.

# Solutions

There is a few options to approach this problem, and this document goes through some of them, from the one which could come first to an inexperienced developer's head to a couple of them with important improvements.

Well see that the first approach, wich implements the trivial solution, has am awful performance, while the second and third one performs thousands of times better with a cost of some extra memory use.

## Solution 1

This approach collects all the words in the dictionary and stores them in a list. In order to find the anagrams for a given word, the algorithm needs to sort each of the words in the dictionary to compare them to the sorted given word.

The building of the list is very fast, as no operation involved. However, further searchs are very slow due to the dictionary needs to be comply walked in order to find anagrams.

```python
 97 class Anagrams1(Anagrams):
 98     """
 99     Very poor performance: This approach collects all the words in the
100     dictionary and stores them in a list.
101     In order to find the anagrams for a given word, the algorithm needs
102     to sort each of the words in the dictionary to compare them to the
103     sorted given word.
104     """
105
106     def __init__(self, source):
107         Anagrams.__init__(self, source)
108         self.words = [w[:-2].lower() for w in open(self.source).readlines()]
109
110     @timing
111     def get_anagrams(self, word):
112         anagrams = []
113         word = "".join(c for c in sorted(word.lower()))
114         for w in self.words:
115             if len(w) != len(word):
116                 continue
117             if "".join(c for c in sorted(w)) == word:
118                 anagrams.append(w)
119         return anagrams
120
```

# Solution 2

In this solution, a python dictionary is created in order to store a pair keys - values, where key is the ordered characters representation of each word in the original dictionary and value is a list containing all the words in the original dictionary where their ordered characters representation is the same that the key.

In this case, collecting the words from the original words dictionary is slightly slower and it requires extra memory ( more or less twice, actually ) but the performance later on, getting the anagrams for a given word is much better as only indexing the characters ordered representation of the given word will return all its anagrams.

```python
123  class Anagrams2(Anagrams):
124      """
125      Much better performance: Create a python dictionary where for each
126      original word in the words dictionary, it stores:
127          - key: the original sorted word
128          - value: all the words that once ordered are the same.
129      """
130
131      def __init__(self, source):
132          Anagrams.__init__(self, source)
133          self.words = {}
134          with open(self.source) as words:
135              for word in [w[:-2].lower() for w in words]:
136                  key = "".join(c for c in sorted(word))
137                  self.words.setdefault(key, [])
138                  self.words[key].append(word)
139
140      @timing
141      def get_anagrams(self, word):
142          key = "".join(c for c in sorted(word.lower()))
143          return self.words.get(key, [])
144
```

## Solution 3

Similarly to solution 2, buils a python dictionary where the key is the hash of the ordered characters representation for each of the original words and value is a list containing all words where the hash of their ordered characters representation matches the key.

This one should be the best approach in performance and the extra memory used for the keys is fixed to *size of integer* * number of words.

```
147 class Anagrams3(Anagrams):
148     """
149     Hash keys: Create a python dictionary where for each
150     original word in the words dictionary, it stores:
151       - key: the hash of the original sorted word
152       - value: all the words that once ordered have the same hash.
153     """
154
155     def __init__(self, source):
156         Anagrams.__init__(self, source)
157         self.hashes = {}
158         with open(self.source) as words:
159             for word in [w[:-2].lower() for w in words]:
160                 key = hash("".join(c for c in sorted(word)))
161                 self.hashes.setdefault(key, [])
162                 self.hashes[key].append(word)
163
164     @timing
165     def get_anagrams(self, word):
166         key = hash("".join(c for c in sorted(word.lower())))
167         return self.hashes.get(key, [])
168
```

# Results

Solution 1, as expected, has a very bad performance.

Running each of the aproaches 500 times, Solution 1 is between 5000 and 8000 times slower than Solution 2 and Solution 3

| ta/tb | Solution 1 | Solution 2 | Solution 3 |
|---|---|---|---|
| Solution1 | | 7763.218794 | 7645.291891 |
| Solution2 | | | 0.984810 |

Solution 2 and Solution 3 are almonst the same, being Solution 2 slightly faster than Solution 3 ( probably because of the cost of hash ).

Solution 3 is, however, less memory consumming.

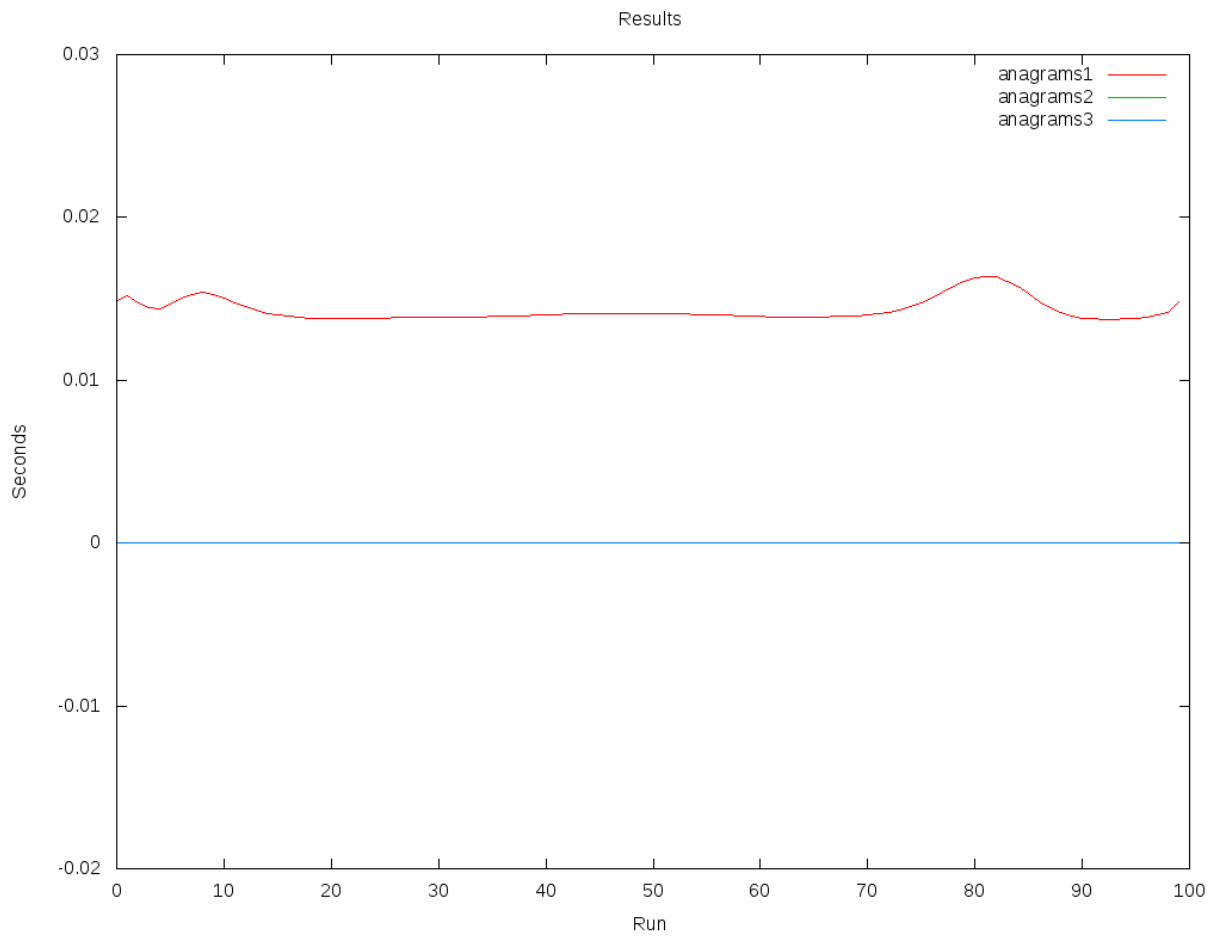Figure 1 representes the times for the three solutions.



*Fig. 1: 50 ran times, solutions 1, 2 and 3*

Figure 2 represents times for solutions 2 and 3. Both solutions present a very similar performance.
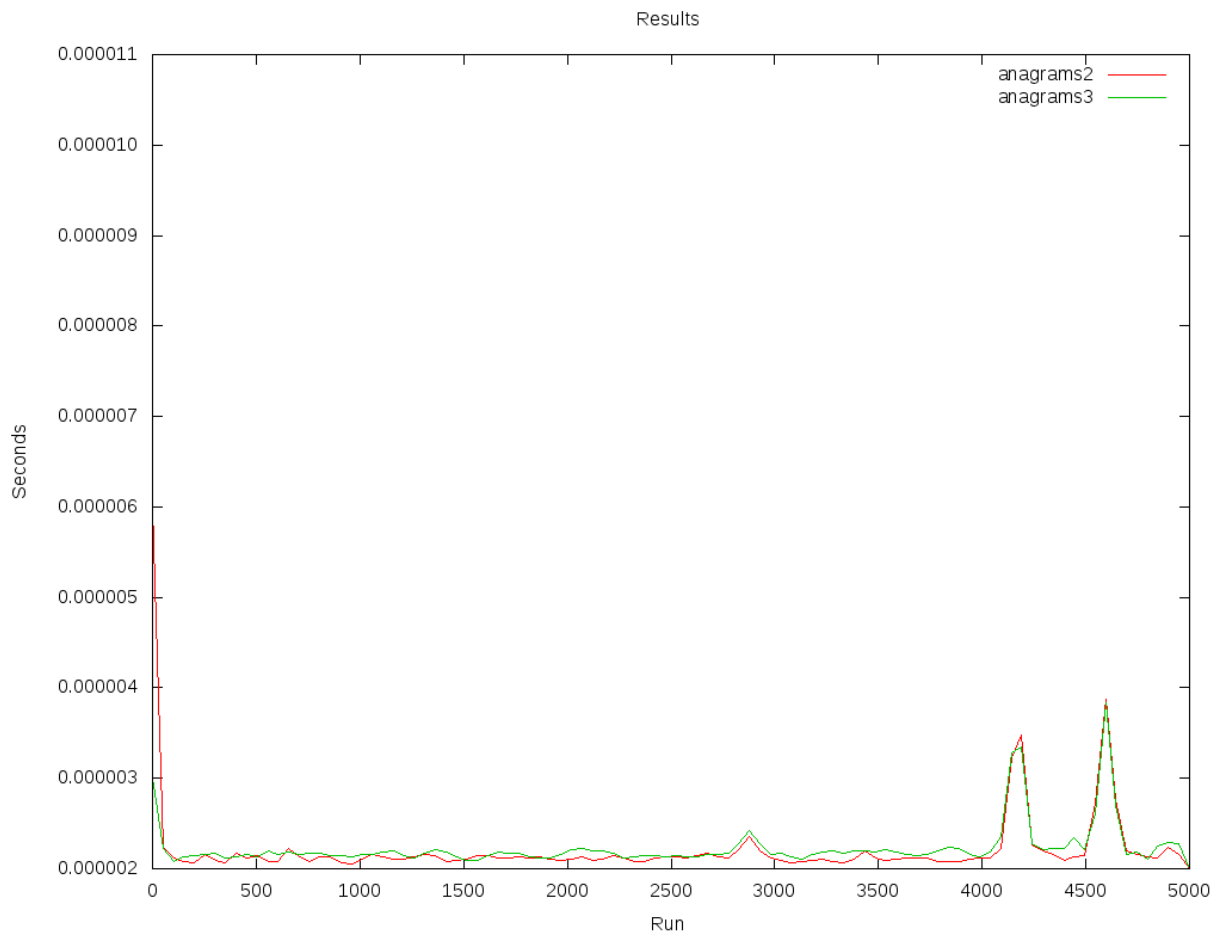


*Fig. 2: 5000 ran times, solutions 2 and 3*

Having a look to these results, the election of Solution 2 or Solution 3 would depend on which is more important in a real proyect:

- Is it critical to be as fast as possible and to use more memory is not a big deal ?

  Solution 2 wins.

- Is it critical to save memory and having a slighty slower algorithm is suitable ?

  Solution 3 wins.

# Latest considerations

- About tests

  An exhaustive test is run covering 100% of the words in the given dictionary

- About threading

  All solutions are thread safe

- About performance

  Solutions 2 and 3 have a very good performance.

# Test environment

- Intel(R) Core(TM) i5-5300U CPU @ 2.30GHz.
- Linux Mint 17
- Python 2.7.6

# Appendix

Complete code is bellow:

Note: exhaustive tests are massive so they are in a separated file exhaustive.py

```python
1 #!/usr/bin/python
2
3 """
4 Given a words.txt file containing a newline-delimited list of dictionary
5 words, please implement the Anagrams class so that the get_anagrams() method
6 returns all anagrams from words.txt for a given word.
7
8 **Bonus requirements:**
9
10   - Optimise the code for fast retrieval
11   - Write more tests
12   - Thread safe implementation
13 """
14
15 import os
16 import time
17 import unittest
18
19
20 CUR_DIR = os.path.dirname(os.path.realpath(__file__))
21 OUTPUT_DIR = os.path.join(CUR_DIR, "../output")
22
23 class PureVirtualMethod(Exception):
24     pass
25
26 def timing(f):
27     def inner(self, *args, **kwargs):
28         t0 = time.time()
29         result = f(self, *args, **kwargs)
30         t = time.time() - t0
31         #print "%s.%s: %0.12f" % (self.__class__.__name__, f.__name__, t)
32         return t, result
33     return inner
34
35 class Statistics(object):
36     """
37     This class implements different statistics for the different solutions
38     It also generates some csv files to be able to process them later on,
39     for example with gnuplot
40     """
41     source = os.path.join(CUR_DIR, 'words.txt')
42
43     def __init__(self):
44         self.anagrams1 = Anagrams1(self.source)
45         self.anagrams2 = Anagrams2(self.source)
46         self.anagrams3 = Anagrams3(self.source)
47         self.workers = [self.anagrams1, self.anagrams2, self.anagrams3]
48
49     def ratios(self):
50         averages = []
```

```python
51          for worker in self.workers:
52              elapsed = 0
53              for i in xrange(500):
54                  t, _ = worker.get_anagrams('plates')
55                  elapsed += t
56              averages.append(elapsed/100.0)
57
58          rat1_2 = averages[0]/averages[1]
59          rat1_3 = averages[0]/averages[2]
60          rat2_3 = averages[1]/averages[2]
61
62          print "1 vs 2: %f" % rat1_2
63          print "1 vs 3: %f" % rat1_3
64          print "2 vs 3: %f" % rat2_3
65
66      def gen_csv_all(self):
67          output_file = os.path.join(OUTPUT_DIR, "anagrams1.csv")
68          output = open(output_file, 'w')
69          output.write("anagrams1,anagrams2,anagrams3\n")
70          for i in xrange(100):
71              t0, _ = self.anagrams1.get_anagrams('plates')
72              t1, _ = self.anagrams2.get_anagrams('plates')
73              t2, _ = self.anagrams3.get_anagrams('plates')
74              output.write("%f, %f, %f\n" %(t0, t1, t2))
75          output.close()
76
77      def gen_csv_best(self):
78          output_file = os.path.join(OUTPUT_DIR, "anagrams2.csv")
79          output = open(output_file, 'w')
80          output.write("anagrams2,anagrams3\n")
81          for i in xrange(5000):
82              t1, _ = self.anagrams2.get_anagrams('plates')
83              t2, _ = self.anagrams3.get_anagrams('plates')
84              output.write("%f, %f\n" %(t1, t2))
85          output.close()
86
87
88 class Anagrams(object):
89
90      def __init__(self, source):
91          self.source = source
92
93      def get_anagrams(self, word):
94          raise PureVirtualMethod("Pure virtual method. Must be redefined")
95
96 # rst-Anagrams1
97 class Anagrams1(Anagrams):
98      """
99      Very poor performance: This approach collects all the words in the
100     dictionary and stores them in a list.
101     In order to find the anagrams for a given word, the algorithm needs
102     to sort each of the words in the dictionary to compare them to the
103     sorted given word.
104     """
105
106     def __init__(self, source):
```

```python
107          Anagrams.__init__(self, source)
108          self.words = [w[:-2].lower() for w in open(self.source).readlines()]
109
110      @timing
111      def get_anagrams(self, word):
112          anagrams = []
113          word = "".join(c for c in sorted(word.lower()))
114          for w in self.words:
115              if len(w) != len(word):
116                  continue
117              if "".join(c for c in sorted(w)) == word:
118                  anagrams.append(w)
119          return anagrams
120
121
122 # rst-Anagrams2
123 class Anagrams2(Anagrams):
124      """
125      Much better performance: Create a python dictionary where for each
126      original word in the words dictionary, it stores:
127          - key: the original sorted word
128          - value: all the words that once ordered are the same.
129      """
130
131      def __init__(self, source):
132          Anagrams.__init__(self, source)
133          self.words = {}
134          with open(self.source) as words:
135              for word in [w[:-2].lower() for w in words]:
136                  key = "".join(c for c in sorted(word))
137                  self.words.setdefault(key, [])
138                  self.words[key].append(word)
139
140      @timing
141      def get_anagrams(self, word):
142          key = "".join(c for c in sorted(word.lower()))
143          return self.words.get(key, [])
144
145
146 # rst-Anagrams3
147 class Anagrams3(Anagrams):
148      """
149      Hash keys: Create a python dictionary where for each
150      original word in the words dictionary, it stores:
151        - key: the hash of the original sorted word
152        - value: all the words that once ordered have the same hash.
153      """
154
155      def __init__(self, source):
156          Anagrams.__init__(self, source)
157          self.hashes = {}
158          with open(self.source) as words:
159              for word in [w[:-2].lower() for w in words]:
160                  key = hash("".join(c for c in sorted(word)))
161                  self.hashes.setdefault(key, [])
162                  self.hashes[key].append(word)
```

```python
163
164        @timing
165        def get_anagrams(self, word):
166            key = hash("".join(c for c in sorted(word.lower())))
167            return self.hashes.get(key, [])
168
169 # rst-Tests
170 class TestAnagrams(unittest.TestCase):
171
172        source = os.path.join(CUR_DIR, 'words.txt')
173
174        def setUp(self):
175            self.anagrams1 = Anagrams1(self.source)
176            self.anagrams2 = Anagrams2(self.source)
177            self.anagrams3 = Anagrams3(self.source)
178
179        def test_no_word(self):
180            _, r =  self.anagrams1.get_anagrams("")
181            self.assertEqual([], r)
182            _, r =  self.anagrams2.get_anagrams("")
183            self.assertEqual([], r)
184            _, r =  self.anagrams3.get_anagrams("")
185            self.assertEqual([], r)
186
187        def test_pure_virtual(self):
188            class AnagramsX(Anagrams):
189                def __init__(self, source):
190                    Anagrams.__init__(self, source)
191
192            anagrams = AnagramsX(self.source)
193            self.assertRaises(PureVirtualMethod, anagrams.get_anagrams, 'pastel')
194
195        def test_exaustive(self):
196            """
197            This tests tests all anagrams for all words in the given
198            dictionary !!
199            They are not ran as individual tests because the required resources
200            for it would be massive.
201            """
202            import exhaustive
203            exhaustive.test_exhaustive(self)
204
205
206 if __name__ == '__main__':
207     unittest.main()
208     #statistics = Statistics()
209     #statistics.ratios()
210     #statistics.gen_csv_all()
211     #statistics.gen_csv_best()
```