

Lab 5 : Computing Similarities Trough Hash Tables and Binary Search Trees

Steven Robles *

University of Texas, El Paso, TX

April 4, 2019

This paper aims to demonstrate implementations of a binary tree and a hash table and compare the performances of both. The data being stored in both data structures is parsed from a huge file. From this information, the data needs to be requested and computed, all of which are timed to evaluate the algorithm's performance.

Contents

I	Introduction	1
II	Design	2
	A Part 1: Binary Search Tree	2
	B Parts 2: Hash Table	2
III	Results	2
IV	Discussion	3
V	Source Code	4
VI	Academic Dishonesty	9

I. Introduction

Binary search trees and hash tables make up the most essential parts of data structures. They are both designed to sort data while they being stored within the structure. Of coarse, they performance will vary depending on the layout of the algorithm and the size of the file. The lab exits to demonstrate and evaluate both data structures while parsing, handling, and computing large amounts of data.

To deploy an accurate analysis of both structures, every function and tasks performed by each data structure is timed and recorded for later comparison. Additional statistics pertaining to the each individual structure such as tree height and load factor will be calculated in order to get a full grasp on its preforming. The computation of the similarities of two data points (words in this case) will also be displayed which shows how accurate the function is able to retrieve the information under a specific certain time.

*ID Number: 80678755

II. Design

This lab is based upon a study that calculated the similarities between two words based upon their numeric embeddings. The scale ranged from -1 to 1 (least similar to most similar). Since the initial file which contains dictionary of words utilized for this lab was large, an initial parsing operation must be preformed. This took some challenge as the file contains 400,000 words with 50 embeddings each. This operations took $O(n^2)$ time.

A. Part 1: Binary Search Tree

With the file already already parsed and its contents stored in an array, a binary search tree (BST) can be constructed to store the data for easy retrieval upon request. One of the challenges that appeared when implementing this structure is how the data is going to be stored within the tree. Since the program is based on python, the final solution was set on comparing the values of the word characters as strings. This process took $O(\log n)$ time.

Additional functions such as calculating the amount of nodes created and the height of the tree also increased the over all time performance of the program. Each of them contributed $O(\log n)$ time as well, although more efficient solutions could have been implemented in their place.

B. Parts 2: Hash Table

The hash table data structure was harder to implement compared to that of the binary search tree. The hash table initially starts at a certain size. In this lab, the size of the hash table starts at 991 *buckets*. It doubles in size every time the load factor reaches 1 but since the hash table is handling a large set of data, it is very inefficient to calculate the load factor of the table every time a new item is inserted. This problem was overcome by tracking how many items are inserted into the table. Clearly not every bucket might be occupied when the lists double in size, but it gives an almost accurate solution which doesn't involve increasing the time execution.

The hash table receives an item, which consists of a word attached to its 50 embeddings, and stores the item in the hash table according to the value of its last character. The average time complexity for this function is $O(n)$, and may increase when it doubles in size. Additional challenges that appeared when implementing this structure was calculating its standard deviation according to its buckets. A formula needed to be created which was able to effectively count all of the times within a bucket and disregard those with no items at all.

III. Results

After the implementations of both data structure types were successfully, each function was executed and recorded for time performance. The time values were used to evaluate each program and compare each data structure with one another. Below are their time executions and overall statistical results.

Figure 1: BST Statistics

```
Binary Search Tree Stats:
Number of Nodes: 400000
Height of Tree : 53
Running Time For Binary Search Tree Construction : 7.377795627020532 Seconds
```

Figure 2: BST Time Performances

```
Average Running Time For Binary Search Tree Query Processing : 2.324043113427857
5e-05 Seconds
Total Time : 0.0006972129340283573 Seconds
```

Figure 3: Hash Table Statistics

```
Hash Table Stats :  
Initial Table Size : 991  
Final Table Size : 507392  
Load Factor : 0.7883451059535822  
Percentage of Empty Lists : 21.165489404641775  
Standard Deviation of the Lengths of the Lists: 18818.258843421572  
Number of Items: 400000
```

Figure 4: Hash Table Time Performances

```
Average Running Time For Hash Table Query Processing : 2.1644867956638336e-05 Seconds  
Total Time : 0.0006493460386991501 Seconds  
Average Search Time : 0.00026499936357140543 Seconds  
Total Time : 0.007949980907142162 Seconds
```

Figure 5: Similarity Computation Outputs From Both Programs

```
Similarity ['bear', 'bear'] : 1.0  
Similarity ['barley', 'shrimp'] : 0.5353  
Similarity ['barley', 'oat'] : 0.6696  
Similarity ['federer', 'baseball'] : 0.287  
Similarity ['federer', 'tennis'] : 0.7168  
Similarity ['harvard', 'stanford'] : 0.8466  
Similarity ['harvard', 'utep'] : 0.0684  
Similarity ['harvard', 'ant'] : -0.0267  
Similarity ['raven', 'crow'] : 0.615  
Similarity ['raven', 'whale'] : 0.3291  
Similarity ['spain', 'france'] : 0.7909  
Similarity ['spain', 'mexico'] : 0.7514  
Similarity ['mexico', 'france'] : 0.5478  
Similarity ['mexico', 'guatemala'] : 0.8114  
Similarity ['computer', 'platypus'] : -0.1277  
Similarity ['computer', 'ant'] : 0.1001
```

IV. Discussion

Overall, implementing both a binary tree and a hash table was a success and was able to handle, store, and retrieve data from the large file provided. The hash table did outperform the binary search tree in terms of executing in the least time, but both implementations ended in the same exact results when computing similarities.

V. Source Code

```
1 """
2 Author: Steven J. Robles
3 Class: CS 2302 Data Structures III
4 Instructor: Olac Fuentes
5 TA: Anindita Nath And Maliheh Zargaran
6 Last Modified: 04/01/2019
7 Discreption: Lab 5:
8     The purpse of this program is to act as the main file of two other porgrams that
9     build a hash table and
10     a binary tree. Each function is called and excecuted upon prompt request
11 """
12
13 from Hashtable import hashTable
14 from B_Tree import treeBuilder
15 from Parse import parseText
16
17 loop = True
18
19 #the following parses the file before hand
20 arr1, arr2 = parseText()
21
22 #the is the while loop which pompts the user.
23 while loop:
24
25     print("Which would rather choose?")
26     print("1. Binary Search Tree")
27     print("2. Hash Table")
28     number = input("3. Exit\n")
29
30     #trys converting the input into an int. if it fails, the pormpt runs again
31     try:
32         choice = int(number)
33     except:
34         choice = -1
35
36     #has table runs
37     if choice == 1:
38         treeBuilder(arr1, arr2)
39
40     #binary tree runs
41     elif choice == 2:
42         hashTable(arr1, arr2)
43
44     #program exits
45     elif choice == 3:
46         print("Good Bye!")
47         loop = False
48     else:
49         print("Try Again")
50
51 """
52 Author: Steven J. Robles
53 Class: CS 2302 Data Structures III
54 Instructor: Olac Fuentes
55 TA: Anindita Nath And Maliheh Zargaran
56 Last Modified: 04/01/2019
57 Discreption: Lab 5:
58     The purpose of this program is to construct a binary search tree based on the parse
59     file from
60     the text. This binary tree will then search for key words anc compute its
61     similarties. It's
62     preformance will be times and compared to that of the hash table
63 """
64
65 from Parse import parseText
```

```

15 import math
16 import timeit
17
18 class BST(object):
19     # Constructor
20     def __init__(self, item, LIn, left=None, right=None):
21         self.item = item
22         self.L = LIn
23         self.left = left
24         self.right = right
25
26 #inserts the item into the bst according to alphabetical order
27 def Insert(T,newItem, listIn):
28     if T == None:
29         T = BST(newItem, listIn)
30     elif T.item > newItem:
31         T.left = Insert(T.left, newItem, listIn)
32     else:
33         T.right = Insert(T.right, newItem, listIn)
34     return T
35
36 #returns the number of items in a BST
37 def CountItems (T):
38     if T is not None:
39         count = 1
40         count += CountItems(T.right)
41         count += CountItems(T.left)
42     return count
43     return 0
44
45 #searches key word within the tree
46 def Search(T, k):
47     while T is not None:
48         if T.item == k:
49             return T
50         if T.item < k:
51             T = T.right
52         elif T.item > k :
53             T = T.left
54     return T
55
56 #computes the similarity between two words using its embeddings
57 def sim(w0, w1):
58     dotSum, sumMag0, sumMag1 = 0.0, 0.0, 0.0
59     for i in range(50):
60         dotSum += w0.L[i] * w1.L[i]
61         sumMag0 += w0.L[i] * w0.L[i]
62         sumMag1 += w1.L[i] * w1.L[i]
63     sumMag0 = math.sqrt(sumMag0)
64     sumMag1 = math.sqrt(sumMag1)
65     return round(dotSum/(sumMag0*sumMag1), 4)
66
67 #computes the height of the tree
68 def Height(T):
69     if T is None:
70         return 0
71
72     else:
73         leftH = Height(T.left)
74         rightH = Height(T.right)
75         if (leftH > rightH): #returns only the biggest count from the left and right sub
76             trees
77             return leftH+1
78         else:
79             return rightH+1
80
81 #main function called from the program
82 def treeBuilder(arr, samples):
83     T = None
84     embed = []

```

```

84     simCount = []
85     time = []
86     nodeCount = 0
87
88     #constructs the tree by inserting nodes
89     start = timeit.default_timer()
90     for i in range(len(arr[0][0])):
91         T = Insert(T, arr[0][0][i], arr[1][0][i])
92         nodeCount +=1
93     stop = timeit.default_timer()
94     constructionT = stop - start
95
96     #retrieves the embeddings of the key words by searching the binary tree
97     for i in range(len(samples)):
98         embed.append([Search(T, samples[i][0]), Search(T, samples[i][1])])
99
100    #computes the similarity between the words by using its embedding information
101    for i in range(len(embed)):
102        start = timeit.default_timer()
103        simCount.append(sim(embed[i][0], embed[i][1]))
104        stop = timeit.default_timer()
105        time.append(stop-start)
106
107    #from here one, the results are printed
108    print("Binary Search Tree Stats: ")
109    print("Number of Nodes: ", nodeCount)
110    print("Height of Tree : ", Height(T))
111    print("Running Time For Binary Search Tree Construciton : ", constructionT, " Seconds")
112    print()
113    print("Reading Word Files to Determine Similarities")
114    print()
115
116    for i in range(len(simCount)):
117        print("Similarity ", samples[i], " : ", simCount[i])
118    print()
119
120    sums = 0
121    for i in range(len(time)):
122        sums += time[i]
123    avg = sums/len(time)
124
125    print("Average Runnig Time For Binary Search Tree Qeury Processing : ", avg, " Seconds")
126    print("Total Time : ", sums, " Seconds")
127    print()

```

```

1  """
2  Author: Steven J. Robles
3  Class: CS 2302 Data Structures III
4  Instructor: Olac Fuentes
5  TA: Anindita Nath And Maliheh Zargaran
6  Last Modified: 04/01/2019
7  Discreption: Lab 5:
8      The purpose of this program is to construct a hash table which stores the information
9      of the
10     parsed file. The size of the table doubles by size each time the laod factor reaches
11     one.
12     The preformance of this function is kept on record to compare to that of the binary
13     search
14     tree implemitation.
15 """
16
17 from Parse import parseText
18 import math
19 import timeit
20
21 class HashTableC(object):
22     def __init__(self, slots):
23         self.item = []
24         self.num_items = 0
25         self.size = slots

```

```

23         for i in range(slots):
24             self.item.append([])
25
26 #resizes the hash table by doubling its current size if the load factor reaches 1
27 def reSize(H):
28     if H.num_items == H.size:
29         for i in range(H.size):
30             H.item.append([])
31         H.item.append([])
32         H.size *=2
33
34 #inserts element onto hash table
35 def InsertC(H,k,l):
36     reSize(H)
37     b = h(k,len(H.item))
38     H.item[b].append([k,l])
39     H.num_items +=1
40
41 #finds element within the hash table and returns its embeddings
42 def FindC(H,k):
43     b = h(k,len(H.item))
44     for i in range(len(H.item[b])):
45         if H.item[b][i][0] == k:
46             return H.item[b][i][:]
47     return b, -1, -1
48
49 #computes the bucket location in which the item will be stored in
50 def h(s,n):
51     r = 0
52     for c in s:
53         r = (r*n + ord(c))% n
54     return r
55
56 #computes the similarity between two words based on their embeddings
57 def sim(w0, w1):
58     dotSum, sumMag0, sumMag1 = 0.0, 0.0, 0.0
59     for i in range(50):
60         dotSum += w0[1][i] * w1[1][i]
61         sumMag0 += w0[1][i] * w0[1][i]
62         sumMag1 += w1[1][i] * w1[1][i]
63     sumMag0 = math.sqrt(sumMag0)
64     sumMag1 = math.sqrt(sumMag1)
65     return round(dotSum/(sumMag0*sumMag1), 4)
66
67 #computes the sandard deviation of the buckets of the hash table
68 def standardDev(H):
69     mean = 0
70     for i in range(len(H.item)):
71         if H.item[i] is not None:
72             mean += len(H.item[i])
73     mean = mean /len(H.item)
74     dist = 0
75     for i in range(len(H.item)):
76         if H.item[i] is not None:
77             dist += abs(len(H.item[i]) - mean) * abs(len(H.item[i]) - mean)
78     return dist/len(H.item)
79
80 #function that is called from the main program
81 def hashTable(arr, samples):
82
83     embed = []
84     simCount = []
85     time = []
86     searchTime = []
87     H = HashTableC(991)
88
89     #fills the hash table with the parsed items from the text file
90     for i in range(len(arr[0][0])):
91         InsertC(H, arr[0][0][i], arr[1][0][i])
92

```

```

93 #finds the embeddings of the key words that are going to be compared
94 for i in range(len(samples)):
95     start = timeit.default_timer()
96     embed.append([FindC(H, samples[i][0]), FindC(H, samples[i][1])])
97     stop = timeit.default_timer()
98     searchTime.append(stop-start)
99
100 #calculates the similarity and times the time it takes
101 for i in range(len(embed)):
102     start = timeit.default_timer()
103     simCount.append(sim(embed[i][0], embed[i][1]))
104     stop = timeit.default_timer()
105     time.append(stop-start)
106
107 #from here on out, the results are displayed
108 print("Hash Table Stats : ")
109 print("Initial Table Size : 991")
110 print("Final Table Size : ", H.size)
111 print("Load Factor : ", H.num_items/H.size)
112 print("Percentage of Empty Lists : ", ((H.size -H.num_items)/H.size)*100)
113 print("Standard Deviation of the Lengths of the Lists: ", standardDev(H))
114 print()
115
116 print("Reading Word File to Determine Similarities")
117 print()
118
119 for i in range(len(simCount)):
120     print("Similarity ", samples[i] , " : ", simCount[i])
121 print()
122
123 sums1, sums2 = 0, 0
124 for i in range(len(time)):
125     sums1 += time[i]
126 avg = sums1/len(time)
127
128 for i in range(len(searchTime)):
129     sums2 += searchTime[i]
130 avg2 = sums2/len(searchTime)
131
132 print("Average Runnig Time For Hash Table Qeury Processing : ", avg, " Seconds")
133 print("Total Time : ", sums1 , " Seconds")
134 print()
135 print("Average Search Time : ", avg2, " Seconds")
136 print("Total Time : ", sums2 , " Seconds")
137 print()

```

```

1 def parseText():
2     f = open('glove.6B.50d.txt',encoding='utf-8')
3
4     lines = f.readlines()
5     for line in f:
6         lines = f.readlines()
7
8     glue = [[[ for x in range(2)]]
9
10    for i in range(len(lines)):
11        test = lines[i].split()
12        glue[0][0].append(test[0])
13        for j in range(50):
14            test[j+1] = float(test[j+1])
15        glue[1][0].append(test[1:])
16
17
18
19 w = open('words.txt',encoding='utf-8')
20
21 lines = w.readlines()
22 for line in w:
23     lines = w.readlines()
24

```



```
25 samples = []
26
27 for i in range(len(lines)):
28     test = lines[i].split()
29     samples.append(test)
30
31 return glue, samples
```

VI. Academic Dishonesty

Scholastic Dishonesty

Any student who commits an act of scholastic dishonesty is subject to discipline. Scholastic dishonesty includes, but not limited to cheating, plagiarism, collusion, the submission for credit of any work or materials that are attributable to another person.

- **Cheating**
 - Copying from the test paper of another student
 - Communicating with another student during a test
 - Giving or seeking aid from another student during a test
 - Possession and/or use of unauthorized materials during tests (i.e. Crib notes, class notes, books, etc)
 - Substituting for another person to take a test
 - Falsifying research data, reports, academic work offered for credit
- **Plagiarism**
 - Using someone's work in your assignments without the proper citations
 - Submitting the same paper or assignment from a different course, without direct permission of instructors
- **Collusion**
 - Unauthorized collaboration with another person in preparing academic assignments

Sign:  Date: 04/04/2019