# CS 2302 Data Structures
# Fall 2019

# Lab Report #5

Due: November 1$^{st}$, 2019
Professor: Olac Fuentes
TA: Anindita Nath

# Introduction

For this lab, we were asked to compare the running times of a hash table with chaining (HTC) and a hash table with linear probing (HTLP), designed to hold an organized list of lists, each of those inner lists containing objects. The objects in question for this lab are those of the WordEmbedding class, in which a word and its vector description, or embedding, is stored. To implement these data structures, we will read a set of 400,000 words and their embeddings from a file and construct instances of these hash tables using the classes and methods defined previously in class and written in previous assignments. The main objective of this lab is to recognize how the nature of each data structure influences its respective running time, as well as to ensure an understanding of how to access any particular element of these data structures.

# Proposed Solution Design and Implementation

To preface, everything is exactly the same as Lab 4 in terms of the initial user prompt, the construction of the desired data structures. All methods necessary for the two classes were already provided, save the *h(k)* method that makes up Part 1.

**Part 1(*h(k)*):**

Of note this time around was that since the program was more or less written up in the last lab, the only compelling part of the story in designing this program was building the *h(k)* method that would determine the hash function by which each variable was stored. Even then, it was simply a matter of writing down the variables as was instructed in the handout.

To elaborate briefly, there were six hash functions to write, all ended with "% the length of the hash table:" The first was simply the length of the inputted word; The second was the ascii value of the word's first character; The third was the product of the ascii values of the word's first and last characters; The fourth was the sum of the ascii values of all the word's letters; the fifth was a recursive function that took the ascii value of the first letter and added the ascii value of each following letter multiplied by 255; and lastly, since the sixth hash function was mine to code, I decided to follow in the footsteps of the third and fourth by taking the first letter's ascii value % the last letter's ascii value.

I did also essentially copy-paste these to the *HashTableLP()* class's *h(k)* method, making sure to change out references to *self.bucket* with *self.item*, but ultimately I was never able to fully construct a HashTableLP because it took so much runtime. I'm not sure if I implemented the class incorrectly, if that's intentional, given part 2, but either way, it was highly impractical to wait and record those runtimes. Suffice it to say, they were much longer and much more inefficient than even HTC(1).

# Experimental Results

**Part 1:**

**h(1)**

```
Choose a hash table implementation.
Type 1 for chaining or 2 for linear probing.

Choice: 1
Choose a hash function corresponding to Part 1 (1-6):

Choice: 1

Hash Table Chain stats:
Running time for hash table chain construction: 517.0052

Running time for hash table chain query processing: 0.8505
```

**h(2)**

```
Choose a hash table implementation.
Type 1 for chaining or 2 for linear probing.

Choice: 1
Choose a hash function corresponding to Part 1 (1-6):

Choice: 2

Hash Table Chain stats:
Running time for hash table chain construction: 234.8463

Running time for hash table chain query processing: 0.3634
```

**h(3)**

```
Choose a hash table implementation.
Type 1 for chaining or 2 for linear probing.

Choice: 1
Choose a hash function corresponding to Part 1 (1-6):

Choice: 3

Hash Table Chain stats:
Running time for hash table chain construction: 76.6083

Running time for hash table chain query processing: 0.0834
```

**h(4)**

```
Choose a hash table implementation.
Type 1 for chaining or 2 for linear probing.

Choice: 1
Choose a hash function corresponding to Part 1 (1-6):

Choice: 4

Hash Table Chain stats:
Running time for hash table chain construction: 35.4035

Running time for hash table chain query processing: 0.0555
```

**h(5)**

```
Choose a hash table implementation.
Type 1 for chaining or 2 for linear probing.

Choice: 1
Choose a hash function corresponding to Part 1 (1-6):

Choice: 5

Hash Table Chain stats:
Running time for hash table chain construction: 224.4353

Running time for hash table chain query processing: 0.3721
```

**h(6)**

```
Choose a hash table implementation.
Type 1 for chaining or 2 for linear probing.

Choice: 1
Choose a hash function corresponding to Part 1 (1-6):

Choice: 6

Hash Table Chain stats:
Running time for hash table chain construction: 174.3625

Running time for hash table chain query processing: 0.2849
```
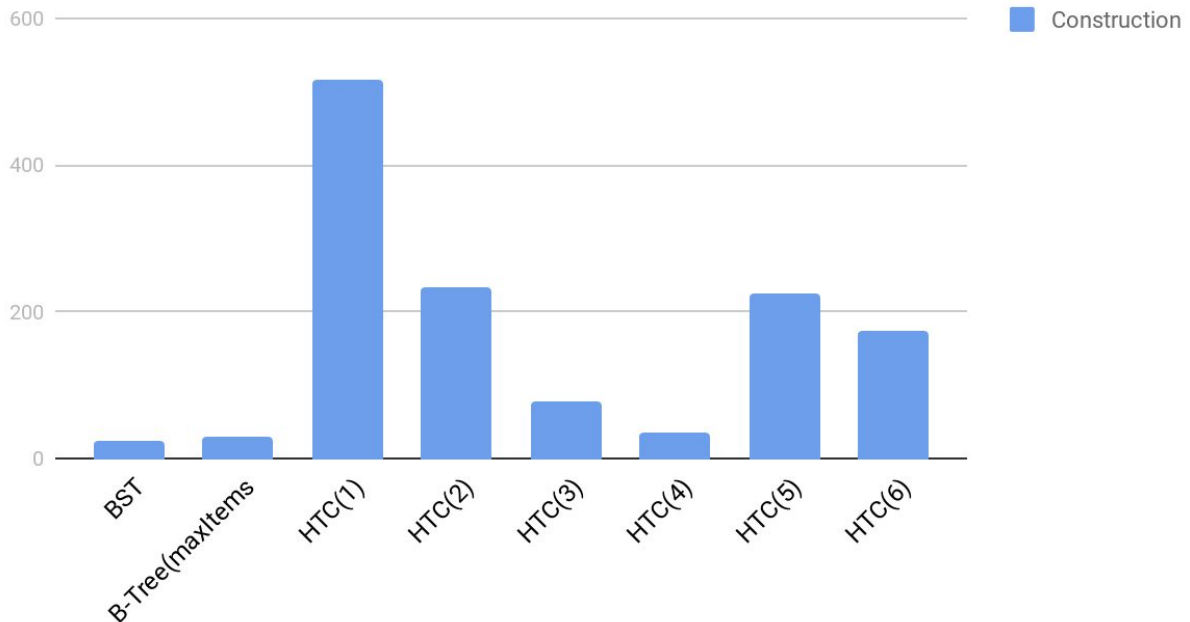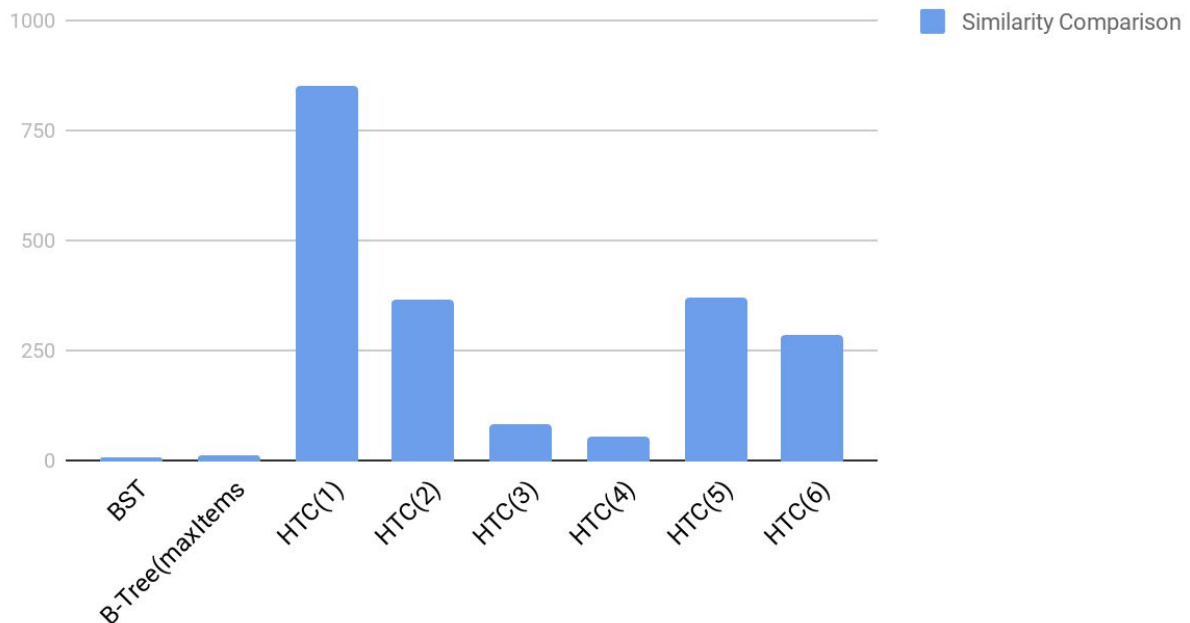
**Tables:**

# Running Times (s)

| Tree/Hash Table | Construction | Similarity Comparison |
| --- | --- | --- |
| BST | 24.4069 | 0.0053 |
| B-Tree(maxItems = 5) | 29.5578 | 0.0133 |
| HTC(1) | 517.0052 | 0.8505 |
| HTC(2) | 234.8462 | 0.36343 |
| HTC(3) | 76.6083 | 0.0834 |
| HTC(4) | 35.0435 | 0.0555 |
| HTC(5) | 224.4353 | 0.3721 |
| HTC(6) | 174.3625 | 0.2849 |

## Runtimes for Construction (s)



## Runtimes for Similarity Comparison (ms)



As the results indicate, every one of the HTC's hash functions failed to live up to either of the trees. HTC(4), the one in which the sum of the ascii values of the letters in a word was

used to sort its insertion, was the only one that came remotely close. Surprisingly, even the recursive function wasn't up to par. Though I'm not exactly what you'd call proud, I am glad at the very least the my own hash function did slightly better than the three worst functions.

# Conclusion

From this lab, I found only trouble regarding the topics. What seemed like an easy enough task at first, as my relative ease with implementing the HTC will demonstrate, quickly become much more difficult to comprehend, as I felt I didn't really understand how to effectively implement the HTLP class. If its poor speed was intentional, I'm still not sold on its viability as a data structure at all, especially when an HTC makes much more sense logically. Perhaps HTLPs are best used with smaller quantities? More than likely, I just failed to properly understand it. That said, I am definitely going to have to study this issue on my own time before too long.

# Appendix

```python
1  # Course: CS 2302 Data Structures
2  # Date of last modification: December 5|
3  # Assignment: Lab 5 Hash Table w/ Chaining/Linear Probing
4  # TA: Anindita Nath
5  # Professor: Olac Fuentes
6  # Purpose: The purpose of this lab is to implement hash tables with chaining and with linear probing in place of
7  #          BSTs and B-trees to compare the running times of all said data structures.
8
9  import numpy as np
10 import time
11
12 #WordEmbedding, HashTableChain, and HashTableLP classes provided by Dr. Fuentes
13 class WordEmbedding(object):
14     def __init__(self,word,embedding):
15         # word must be a string, embedding can be a list or and array of ints or floats
16         self.word = word
17         self.emb = np.array(embedding, dtype=np.float32)
18
19 class HashTableChain(object):
20     # Builds a hash table of size 'size'
21     # Item is a list of (initially empty) lists
22     # Constructor
23     def __init__(self,size, c = 0):
24         self.bucket = [[] for i in range(size)]
25         self.c = c
26
27     def h(self,k):
28         #Part 1
29         #The length of the string % n
30         if self.c == 1:
31             return len(k.word) % len(self.bucket)
32         #The ascii value (ord(c)) of the first character in the string % n
33         elif self.c == 2:
34             return ord(k.word[0]) % len(self.bucket)
35         #The product of the ascii values of the first and last characters in the string % n
36         elif self.c == 3:
37             return (ord(k.word[0]) * ord(k.word[-1])) % len(self.bucket)
38         #The sum of the ascii values of the characters in the string % n
39         elif self.c == 4:
40             ordAll = 0
41             for l in range(len(k.word)):
42                 ordAll += ord(k.word[l])
43             return ordAll % len(self.bucket)
44         #The recursive formulation h(",n) = 1; h(S,n) = (ord(s[0]) + 255*h(s[1:],n))% n
45         elif self.c == 5:
46             if k.word == '':
47                 return 1
48             else:
49                 s = WordEmbedding(k.word[1:], 0)
50                 return (ord(k.word[0]) + 255 * self.h(s)) % len(self.bucket)
51         #Another function of your choice(The product of ascii values of characters in string % n)
52         else:
53             return (ord(k.word[0]) % ord(k.word[-1])) % len(self.bucket)
```

```python
54                                          #to sort WordEmbedding objects alphabetically by ASCII value
55
56      def insert(self,k):
57          # Inserts k in appropriate bucket (list)
58          # Does nothing if k is already in the table
59          b = self.h(k)
60          if not k.word in self.bucket[b]:#Modified "k" to "k.word" for proper detection
61              self.bucket[b].append(k)          #Insert new item at the end
62
63      def find(self,k):
64          # Returns bucket (b) and index (i)
65          # If k is not in table, i == -1
66          b = self.h(k)
67          try:
68              i = self.bucket[b].index(k)
69          except:
70              i = -1
71          return b, i
72
73      def print_table(self):
74          print('Table contents:')
75          for b in self.bucket:
76              print(b)
77
78      def delete(self,k):
79          # Returns k from appropriate list
80          # Does nothing if k is not in the table
81          # Returns 1 in case of a successful deletion, -1 otherwise
82          b = self.h(k)
83          try:
84              self.bucket[b].remove(k)
85              return 1
86          except:
87              return -1
88
89  class HashTableLP(object):
90      # Builds a hash table of size 'size', initilizes items to -1 (which means empty)
91      # Constructor
92      def __init__(self,size,c = 0):
93          self.item = np.zeros(size,dtype=np.int)-1
94          self.c = c
95
96      def insert(self,k):
97          # Inserts k in table unless table is full
98          # Returns the position of k in self, or -1 if k could not be inserted
99          for i in range(len(self.item)): #Despite for loop, running time should be constant for table with low load factor
100             pos = self.h(k)+i
101             numK = []
102             for i in range(len(k.emb)):
103                 numK = ''.join(str(k.emb[i]))
104             numK = float(numK)
105             if self.item[pos] < 0:
106                 self.item[pos] = numK
107                 return pos
```

```python
108            return -1
109
110     def find(self,k):
111         # Returns the position of k in table, or -1 if k is not in the table
112         for i in range(len(self.item)):
113             pos = self.h(k)+i
114             if self.item[pos] == k:
115                 return pos
116             if self.item[pos] == -1:
117                 return -1
118         return -1
119
120     def delete(self,k):
121         # Deletes k from table. It returns the position where k was, or -1 if k was not in the table
122         # Sets table item where k was to -2 (which means deleted)
123         f = self.find(k)
124         if f >=0:
125             self.item[f] = -2
126         return f
127
128     def h(self,k):
129         #Part 1
130         #The length of the string % n
131         if self.c == 1:
132             return len(k.word) % len(self.item)
133         #The ascii value (ord(c)) of the first character in the string % n
134         elif self.c == 2:
135             return ord(k.word[0]) % len(self.item)
136         #The product of the ascii values of the first and last characters in the string % n
137         elif self.c == 3:
138             return (ord(k.word[0]) * ord(k.word[-1])) % len(self.item)
139         #The sum of the ascii values of the characters in the string % n
140         elif self.c == 4:
141             ordAll = 0
142             for l in range(len(k.word)):
143                 ordAll += ord(k.word[l])
144             return ordAll % len(self.item)
145         #The recursive formulation h('',n) = 1; h(S,n) = (ord(s[0]) + 255*h(s[1:],n))% n
146         elif self.c == 5:
147             if k.word == '':
148                 return 1
149             else:
150                 s = WordEmbedding(k.word[1:], 0)
151                 return (ord(k.word[0]) + 255 * self.h(s)) % len(self.item)
152         #Another function of your choice
153         else:
154             return #Modified "k%len(self.bucket)" to "ord(k.word[0])"
155                             #to sort WordEmbedding objects alphabetically by ASCII value
156
157
158     def print_table(self):
159         print('Table contents:')
160         for b in self.bucket:
161             print(b)
```

```python
162
163 #My Code
164 if __name__ == "__main__":
165     #Read file and create table based on tableChoice
166     file = open("/Users/johanncampos369/Downloads/glove.6B/glove.6B.50d.txt", "r")
167
168     while(True):
169         print("Choose a hash table implementation.\nType 1 for chaining or 2 for linear probing.")
170         tableChoice = int(input("Choice: "))
171         print("Choose a hash function corresponding to Part 1 (1-6): ")
172         c = int(input("Choice: "))
173         #Store each WordEmbedding object in appropriate type of table
174         if tableChoice == 1:
175             start = time.time()
176             htc = HashTableChain(255, c)#Set size to 255 to get buckets corresponding to ASCII values
177                                         #Set c to determine which hash function to use.
178             for f in file:
179                 words = f.split(" ")
180                 htc.insert(WordEmbedding(words[0], words[1:]))#Insert new object with word and emb list
181             end = time.time()
182             #Print stats of resulting
183             print("\nHash Table Chain stats: ")
184             print("Running time for hash table chain construction: %.4f" % (end - start))
185             break
186         elif tableChoice == 2:
187             #size = int(input(("Enter size of htlp: ")))
188             print("\nBuilding Hash Table w/ Linear Probing.")
189             start = time.time()
190             htlp = HashTableLP(400001, c)
191             for f in file:
192                 words = f.split(" ")#Split each line into list of word and emb
193                 htlp.insert(WordEmbedding(words[0], words[1:]))#Insert new object with word and emb list
194             end = time.time()
195             #Print stats of resulting B-tree
196             print("\nHash Table LP Stats: ")
197             print("Running time for hash table LP construction: %.4f" % (end - start))
198             break
199         else:
200             print("Invalid entry. Please try again.")
201     #Read second file containing pairs of words
202     print("\nReading word file to determine similarities")
203     file2 = open("/Users/johanncampos369/Downloads/glove.6B/wordpairs.txt", "r")
204     #Compare words in file2 line by line, finding them in tree made from first file
205     print("\nWord similarities found: ")
206     total = 0.0
207     lines = file2.read().splitlines()#Make list of lines without \n
208     for i in range(len(lines)):
209         start = time.time()
210         words = lines[i].split(" ")#Makes list of two words for given line
211         if tableChoice == 1:#htc
212             word1b, word1i = htc.find(WordEmbedding(words[0], 0))#find returns bucket and index of WordEmbedding objects
213             word2b, word2i = htc.find(WordEmbedding(words[1], 0))
214             #Take necessary numbers to calculate similarity based on emb by accessing htc elements via bucket and index
215             top = np.dot(htc.bucket[word1b][word1i].emb, htc.bucket[word2b][word2i].emb)

216             bottom = np.dot(np.linalg.norm(htc.bucket[word1b][word1i].emb), np.linalg.norm(htc.bucket[word2b][word2i].emb))
217         else:#Previous try-except allows for tableChoice to only be either 1 or 2, thus making this htlp
218             word1 = htlp.find(WordEmbedding(words[0],0))#find returns index of WordEmbedding object that contains sought for word
219             word2 = htlp.find(WordEmbedding(words[1],0))
220             #Take necessary numbers to calculate similarity based on emb by accessing htlp elements via index
221             top = np.dot(htlp.item[word1], htlp.item[word2])
222             bottom = np.dot(np.linalg.norm(htlp.item[word1]), np.linalg.norm(htlp.item[word2]))
223         #Calculate similarity based on emb from objects of either table
224         sim = top / bottom
225         end = time.time()
226         total += end - start#Accumulates time taken for every similarity comparison
227         print("Similarity ", words, " = %.4f" % sim)
228     if tableChoice == 1:
229         print("Running time for hash table chain query processing: %.4f" % total)
230     else:
231         print("Running time for hash table LP query processing: %.4f" % total)
```

I certify that this project is entirely my own work. I wrote, debugged, and tested the code being presented, performed the experiments, and wrote the report. I also certify that I did not share my code or report or provided inappropriate assistance to any student in the class.