# CS 2302 Data Structures
# Fall 2019

# Lab Report #4

# Introduction

For this lab we were tasked with creating either a binary search tree or a B-tree depending on what the user whishes to use. After the type of tree was specified, we needed to gather all words without special characters from the file "glove.6B.50b.txt" and create Word-Embedding object to then add into the trees. Once the tree was completed it would output statistics of the tree then we could see how long the tree takes to display the vector differences of two different words as well as the running time needed to do so.

# Proposed Solution Design and Implementation

## Part 1:

First, I created a prompt to ask the user what tree they would like to use as well as ask for the max number of items to store per node if they chose a B-tree.

## Part 2:

Next, I created two more files called "**BTree.py**" and "**BSTree.py**" to create the tree chosen. Within each file there was a class objects that imitated the tree type as well as stored data such as children for a B-tree and left/right for a BST, as well as multiple methods needed to do certain things such as find the height of the tree or the number of nodes within the tree. Then I went line by line in the file "**glove.6B.50b.txt**". I would split the file into two items using the method **lineSplit()**. The first item being returned was a string of the word called "word" and the second was call "emb" which was a float list that contained the words embeddings. After that the word would be checked with the "**.isalpha**()" method to see if it is free of any special characters and if so would be pushed into a Word-Embedding object and placed into the assigned tree type. Once all three hundred thousand plus items were added into the tree it would print out the tree's statistics such as the number of nodes in the tree and the height of the tree. Along with the time taken to create the tree.

## Part 3:

Then, after the tree was created and filled with all the Word-Embedding objects the program would read a new file containing two words per line to compare and print out their similarity. It would split the line into a list of the two words and search through the tree created to find the words embeddings. Next, I would push the two word embedding objects of the words through a method called "**Similarity()**," which would them computer the two words dot product and divided it by the words vector magnitude. It would then print out the two words and their computed similarity. It would repeat this until all line in the file were compared.

# Experimental Results

       For both tree's I ran the program with four test files. The first comparing 15 words, the second comparing 60 words, the third comparing 120 words, and the last comparing 240 words. The first 15 words are the words given to use from our lab instructions and the rest were chosen using a random word generator.

**Binary Search Tree:**

Time Taken to create the Tree: On average 16.1203 Seconds

- `Running time for Binary Search Tree construction: 16.1641  seconds`

**Test 1:**

List Length: 15 words

Time Taken to Compute Similarity's:

- `Running time for Binary Search Tree query processing:  0.0050893`

**Test 2:**

List Length: 60 words

Time Taken to Compute Similarity's:

- `Running time for Binary Search Tree query processing:  0.0060132`

**Test 3:**

List Length: 120 words

Time Taken to Compute Similarity's:

- `Running time for Binary Search Tree query processing:  0.0105908`

**Test 4:**

List Length: 240 words

Time Taken to Compute Similarity':

- `Running time for Binary Search Tree query processing:  0.0238645`

**B-Tree: (Max_items = 5)**

Time Taken to create the Tree: On average 18.6536 Seconds

- `Running time for B-tree construction (with max_items = 5): 18.7628 seconds`

**Test 1:**

List Length: 15 words

Time Taken to Compute Similarity's:

- `Running time for B-Tree query processing (with max_items = 5): 0.0070624 seconds`

**Test 2:**

List Length: 60 words

Time Taken to Compute Similarity's:

- `Running time for B-Tree query processing (with max_items = 5): 0.0109735 seconds`

**Test 3:**

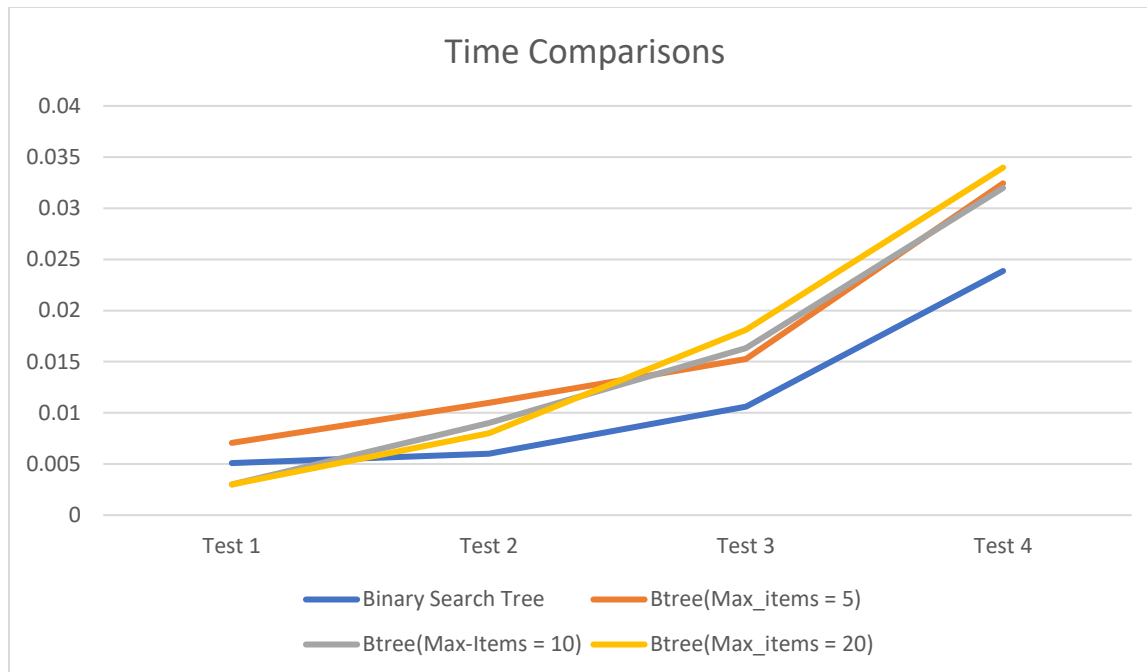List Length: 120 words

Time Taken to Compute Similarity's:

- `Running time for B-Tree query processing (with max_items = 5): 0.0152524 seconds`

**Test 4:**

List Length: 240 words

Time Taken to Compute Similarity's:

- `Running time for B-Tree query processing (with max_items = 5): 0.0324318 seconds`

Time Comparisons

The data collected helps show that the binary search tree always computes data faster then a b-tree with max inputs of 5. But if the Number of items being compared is moderately low and the max items is high then it has a chance of betting the binary search time as well as the b-tree with max input of 5. This is displayed by the two lines for the b-tree with max items 10, and 20.

## Conclusion

In conclusion, I was able to see how in certain instances of searching through trees that b-trees could be faster than binary search trees depending on factors such as the max number of items per node and the number of items being searched for. However, for a more reliable time of completion the binary search tree was the fastest and most consistent.

## Appendix

**Lab4.py:**

```
1    '''
2
3        Cs2302 Data Structures
4        Issac Rivas
5        Lab 4
6        Dr.Fuentues
7
8    '''
9    import BSTree
10   import BTree
11   import numpy as np
12   import time as t
```

```python
13
14   #Creates an object that holds the word as well as 50 element embedding list
15   class WordEmbedding(object):
16       def __init__(self, word, embedding):
17           self.word = word
18           self.emb = np.array(embedding, dtype=np.float32)
19
20   #Splis a text file line into a word and a float list of the remainig 50 elements
21   def lineSplit(line):
22       emb = line.split()
23       word = emb[0]
24       del emb[0]
25       for x in range(0, len(emb)):
26           emb[x] = float(emb[x])
27       return word, emb
28
29   #Calculates for the Similartiy of two words
30   def Similarity(word0, word1):
31       dotProduct = np.dot(word0.emb,word1.emb)
32       magnitude = np.linalg.norm(word0.emb) * np.linalg.norm(word1.emb)
33       return dotProduct/magnitude
34
35   if __name__ == "__main__":
36       print("Choose table implementation:")
37       print("Enter '1' for binary search or '2' for B-tree")
38       ans = int(input("Input: "))
39       while(ans != 1 and ans != 2):
40           print()
41           print("Incorrect input entered please Try again.")
42           print("Type 1 for binary search or 2 for B-tree")
43           ans = int(input("Input: "))
44
45       filepath = 'C:/Users/Issac/Desktop/Lab_4/glove.6B.50d.txt'
46       testpath = 'C:/Users/Issac/Desktop/Lab_4/test.txt'
47       start = t.time()
48
49       #If user wants to create a BST
50       if(ans == 1):
51           T = None
52           with open(filepath,  encoding="utf8") as fp:
53               for line in fp:
54                   word, emb = lineSplit(line)
55                   if word.isalpha():
56                       temp = WordEmbedding(word, emb)
57                       T = BSTree.Insert(T, temp)
58           Total = t.time() - start
59           print()
```

```python
59              print()
60              print("Building Binary Search Tree")
61              print()
62              print("Binary Tree Stats:")
63              print("Number of Nodes: ", BSTree.numOfNodes(T))
64              print("Tree Height is: ", BSTree.treeHeight(T))
65              print("Running time for Binary Search Tree construction:", round(Total, 4), " seconds")
66              print()
67              print("Reading word file to determine similarities")
68              print()
69              print("Word similarities found:")
70              start = t.time()
71              with open(testpath) as tp:
72                  for line in tp:
73                      words = line.split()
74                      temp0 = BSTree.Search(T, words[0])
75                      temp1 = BSTree.Search(T, words[1])
76                      sim = Similarity(temp0.data, temp1.data)
77                      print("Similarity [" + words[0] + "," + words[1] + "] = ", round(sim, 4))
78              Total = t.time() - start
79              print()
80              print("Running time for Binary Search Tree query processing: ", round(Total, 7))
81
82          #If user wants to create a B-tree
83          else:
84              max_data = int(input("Please enter the maximum number of items per node: "))
85              while (max_data < 1):
86                  print("Error please insert a number equal to or greater then 1.")
87                  max_data = int(input("Please enter the maximum number of items per node: "))
88              start = t.time()
89              T = BTree.BTree([], max_data)
90              with open(filepath,  encoding="utf8") as fp:
91                  for line in fp:
92                      word, emb = lineSplit(line)
93                      if word.isalpha():
94                          temp = WordEmbedding(word, emb)
95                          BTree.Insert(T, temp)
96              print()
97              print("Building B-Tree")
98              print()
99              print("B-Tree Stats:")
100             print("Number of Nodes: ", BTree.NumItems(T))
101             print("Tree Height is: ", BTree.Height(T))
102             Total = t.time() - start
103             print("Running time for B-tree construction (with max_items = " + str(max_data) + "):", round(Total, 4)
104             print()
105             print("Reading word file to determine similarities")
```

```
106          print()
107          print("Word similarities found:")
108          start = t.time()
109          with open(testpath) as tp:
110              for line in tp:
111                  words = line.split()
112                  wordemb1 = WordEmbedding(words[0], [])
113                  wordemb2 = WordEmbedding(words[1], [])
114                  temp0 = BTree.Search(T, wordemb1)
115                  temp1 = BTree.Search(T, wordemb2)
116                  sim = Similarity(temp0, temp1)
117                  print("Similarity [" + words[0] + "," + words[1] + "] = ", round(sim, 4))
118          Total = t.time() - start
119          print()
120          print("Running time for B-Tree query processing (with max_items = " + str(max_data) + "):", round(Total,
121
122
123
124
```

**BSTree.py:**

```
1    # -*- coding: utf-8 -*-
2    """
3    Created on Tue Oct 22 19:35:08 2019
4
5    @author: Issac
6    """
7    import matplotlib.pyplot as plt
8    import numpy as np
9
10   class BST(object):
11       def __init__(self, data, left=None, right=None):
12           self.data = data
13           self.left = left
14           self.right = right
15
16   def Insert(T,newItem):
17       if T == None:
18           T =  BST(newItem)
19       elif T.data.word >= newItem.word:
20           T.left = Insert(T.left,newItem)
21       else:
22           T.right = Insert(T.right,newItem)
23       return T
24
```

```python
25    def numOfNodes(T):
26        if T is None:
27            return 0
28        return numOfNodes(T.left) + 1 + numOfNodes(T.right)
29
30    def treeHeight(T):
31        if T != None:
32            left = 1 + treeHeight(T.left)
33            right = 1 + treeHeight(T.right)
34            if left > right:
35                return 1 + left
36            else:
37                return 1 + right
38        else:
39            return -1
40
41    def Search(T, k):
42        if T == None:
43            return T
44        if T.data.word == k:
45            return T
46        if T.data.word > k:
47            return Search(T.left, k)
48        else:
49            return Search(T.right, k)
50
51
52    def DrawBST_(T, x0, x1, y, y_inc,ax):
53        if T is not None:
54            xm = (x0+x1)/2
55            yn = y-y_inc
56            if T.left is not None:
57                p=np.array([[xm,y], [(x0+xm)/2,yn]])
58                ax.plot(p[:,0],p[:,1],linewidth=1,color='k')
59                DrawBST_(T.left,x0,xm,yn, y_inc,ax)
60            if T.right is not None:
61                p=np.array([[xm,y], [(x1+xm)/2,yn]])
62                ax.plot(p[:,0],p[:,1],linewidth=1,color='k')
63                DrawBST_(T.right,xm,x1,yn, y_inc,ax)
64            ax.text(xm,y, str(T.data.word), size=10,ha="center", va="center",
65                bbox=dict(facecolor='w',boxstyle="circle"))
66
67    def DrawBST(T):
68        fig, ax = plt.subplots()
69        DrawBST_(T, 0, 200, 400, 20, ax)
70        ax.set_aspect(1.0)
71        ax.axis('off')
```

```python
72          plt.show()
```

**BTree.py:**

```python
1   # -*- coding: utf-8 -*-
2   """
3   Created on Tue Oct 22 19:18:15 2019
4
5   @author: Issac
6   """
7   class BTree(object):
8       def __init__(self,data,max_data,child=[],isLeaf=True):
9           self.data = data
10          self.child = child
11          self.isLeaf = isLeaf
12          if max_data <3: #max_data must be odd and greater or equal to 3
13              max_data = 3
14          if max_data%2 == 0: #max_data must be odd and greater or equal to 3
15              max_data +=1
16          self.max_data = max_data
17
18  def Insert(T,i):
19      if not IsFull(T):
20          InsertInternal(T,i)
21      else:
22          m, l, r = Split(T)
23          T.data =[m]
24          T.child = [l,r]
25          T.isLeaf = False
26          k = FindChild(T,i)
27          InsertInternal(T.child[k],i)
28
29  def FindChild(T,k):
30      # Determines value of c, such that k must be in subtree T.child[c], if k is in the BTree
31      for i in range(len(T.data)):
32          if k.word < T.data[i].word:
33              return i
34      return len(T.data)
35
36  def FindChildSearch(T,k):
37      # Determines value of c, such that k must be in subtree T.child[c], if k is in the BTree
38      for i in range(len(T.data)):
```

```python
             if k.word < T.data[i].word:
                 return i
         return len(T.data)

     def InsertInternal(T,i):
         # T cannot be Full
         if T.isLeaf:
             InsertLeaf(T,i)
         else:
             k = FindChild(T,i)
             if IsFull(T.child[k]):
                 m, l, r = Split(T.child[k])
                 T.data.insert(k,m)
                 T.child[k] = l
                 T.child.insert(k+1,r)
                 k = FindChild(T,i)
             InsertInternal(T.child[k],i)

     def Split(T):
         #print('Splitting')
         #PrintNode(T)
         mid = T.max_data//2
         if T.isLeaf:
             leftChild = BTree(T.data[:mid],max_data=T.max_data)
             rightChild = BTree(T.data[mid+1:],max_data=T.max_data)
         else:
             leftChild = BTree(T.data[:mid],T.max_data,T.child[:mid+1],T.isLeaf)
             rightChild = BTree(T.data[mid+1:],T.max_data,T.child[mid+1:],T.isLeaf)
         return T.data[mid], leftChild,  rightChild

     def InsertLeaf(T,i):
         temp = []
         for x in range(0, len(T.data)):
             temp.append(T.data[x].word)
         temp.append(i.word)
         temp.sort()
         for y in range(0,len(temp)):
             if temp[y] == i.word:
                 T.data.insert(y, i)
                 return

     def IsFull(T):
         return len(T.data) >= T.max_data

     def Leaves(T):
         # Returns the leaves in a b-tree
         if T.isLeaf:
```

```python
 86            return [T.data]
 87        s = []
 88        for c in T.child:
 89            s = s + Leaves(c)
 90        return s
 91
 92    def Height(T):
 93        if T.isLeaf:
 94            return 0
 95        return 1 + Height(T.child[0])
 96
 97    def PrintD(T,space):
 98        # Prints data and structure of B-tree
 99        if T.isLeaf:
100            for i in range(len(T.data)-1,-1,-1):
101                print(space,T.data[i].word)
102        else:
103            PrintD(T.child[len(T.data)],space+'   ')
104            for i in range(len(T.data)-1,-1,-1):
105                print(space,T.data[i].word)
106                PrintD(T.child[i],space+'   ')
107
108    def Search(T,k):
109        # Returns node where k is, or None if k is not in the tree
110        for x in range(0, len(T.data)):
111            if k.word == T.data[x].word:
112                return T.data[x]
113        if T.isLeaf:
114            return None
115        return Search(T.child[FindChildSearch(T,k)],k)
116
117    def Set_x(T,Dx):
118        # Finds x-coordinate to display each node in the tree
119        if T.isLeaf:
120            return
121        else:
122            for c in T.child:
123                Set_x(c,Dx)
124            d = (Dx[T.child[0].data[0]] + Dx[T.child[-1].data[0]] + 10*len(T.child[-1].data))/2
125            Dx[T.data[0]] = d - 10*len(T.data)/2
126
127    def NumItems(T):
128        if T.isLeaf:
129            return len(T.data)
130        else:
```

```
132            for x in range(0, len(T.child)):
133                totalBelow = totalBelow + NumItems(T.child[x])
134            return totalBelow
```

I Issac Rivas, certify that this project is entirely my own work. I wrote, debugged, and tested the code being presented, preformed 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.