
CS231: Project 6

Word Frequency

Parth Parth | CS231L-A | Dr. Alan Harper

Abstract

In the project, a Generic Binary Search Tree has been created and implemented. A Binary Search Tree is a recursive, node based data structure that stores data according to the following rules: the left subtree has values smaller than the root, the right subtree has values larger than the root, and both subtrees are also binary search trees.

This data structure is specially useful for the use case of this project as it allows us to store and manipulate data in a more efficient way than a list. This data structure has then been used to analyze Reddit comments from years 2008 – 2015. The results have then been tabulated and analyzed.

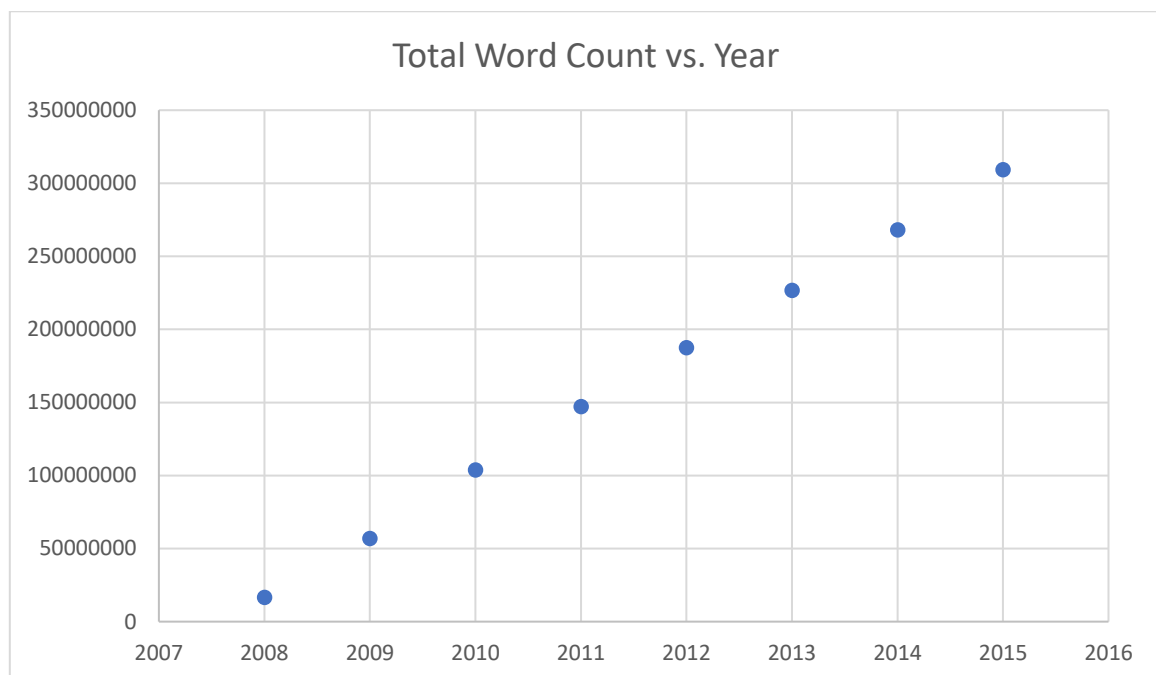
Analysis

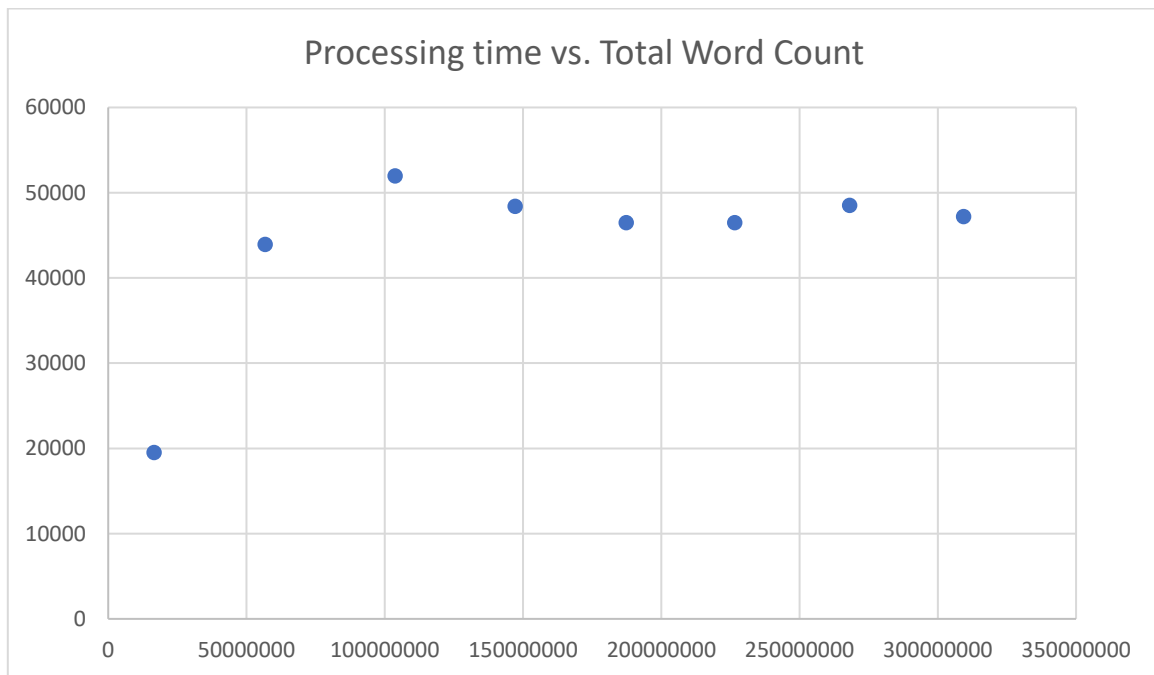
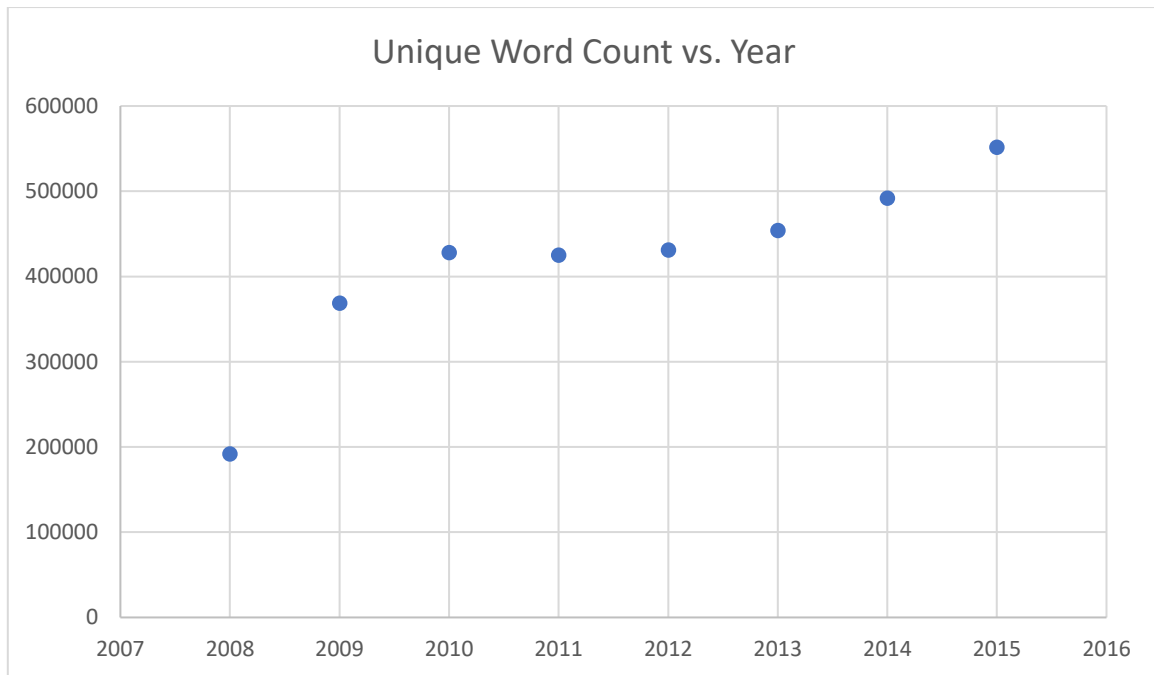
The diff generated from the files obtained from *WCTest.java* is shown in the image below:

```
C:\Users\Parth> FC "C:\Users\Parth\OneDrive\College Classes\CS231\Projects\Project 6\counts_ct.txt" "C:\Users\Parth\OneDrive\College Classes\CS231\Projects\Project 6\counts_ct_v2.txt"
Comparing files C:\USERS\PARTH\ONEDRIVE\COLLEGE CLASSES\CS231\PROJECTS\PROJECT 6\counts_ct.txt and C:\USERS\PARTH\ONEDRIVE\COLLEGE CLASSES\CS231\PROJECTS\PROJECT 6\COUNTS_CT_V2.TXT
FC: no differences encountered
```

So, no differences were found.

The graphs required are as follows:





It is clear from the last graph that the times come closer together as words increase. Since all the words are selected from Reddit, they should be randomly uniformly randomly distributed, giving us a somewhat balanced tree.

The time being recorded is the time taken to read a file (which should be almost constant), add all the words into the tree (which is $O(n \cdot \log(n))$ for the average case). Therefore, it makes sense that the values somewhat look like the graph of $n \cdot \log(n)$.

The reason the graph is not exactly like $n \cdot \log(n)$ might also be because the number of unique words stay almost constant for 2010, 2011, and 2012 where a dip is observed in the third graph.

A separate observation is that the number of total words has grown uniformly, in a linear fashion.

Extension 1

In this extension, I have implemented a *remove(K key)* method that takes in a key of type *K* and removes and returns the node where the key was found if it exists, otherwise it returns null. It uses the following logic:

1. If node has no children, just remove the reference to it in the parent.
2. If it has only one child (or subtree), replace the reference to that node with a reference to the subtree root in the parent
3. If node has both children
 - a. Find minimum in right subtree
 - b. Copy the value to the target node
 - c. Remove the minimum in the right subtree

The following image shows a demonstration:

```
20 Frequency: 20
10 Frequency: 10
5 Frequency: 5
4 Frequency: 4
6 Frequency: 7
9 Frequency: 9
11 Frequency: 11

Size: 7
Removing 6
20 Frequency: 20
10 Frequency: 10
5 Frequency: 5
4 Frequency: 4
9 Frequency: 9
11 Frequency: 11

Size: 6
```

Extension 2

In this extension, I have edited the *toString()* method to return a *String* that contains a breadth-first traversal of the tree i.e. from left to right, level-by-level.

The following image demonstrates this:

```
4 5 6 9 10 11 20
Size: 7
Removing 6
4 5 9 10 11 20
Size: 6
Level-by-level traversal:
Key: 20 Value: 20
Key: 10 Value: 10
Key: 5 Value: 5
Key: 11 Value: 11
Key: 4 Value: 4
Key: 9 Value: 9
```

The first two lines are the inorder traversal of the BST before and after removing the value 6 from the tree.

References/Acknowledgements

The algorithms for the remove method for extension 1 and the level-by-level traversal that I have employed were covered by Dr. Al Madi in class.