

Csci 335 Assignment 2

Due Thursday, March 9

****Programming: Using and Comparing Tree Implementations (100 points)**

The goal of this assignment is to become familiar with trees and compare the performance of the basic binary search tree with the self-balancing AVL tree. You will also work with a real world data set and construct a generic test routine for comparing several different implementations of the tree container class. You are encouraged to use the book's implementation for the BST and AVL trees. You can get the code from:

https://users.cs.fiu.edu/~weiss/dsaa_c++4/code/

(files BinarySearchTree.h, AvlTree.h, dsexceptions.h)

The above files are already included in the code provided to you. Note that you should use the Makefile as provided. If you want to do changes please consult with me first.

Part 1 (10 points)

First, create a class object named SequenceMap that has as private data members the following two:

```
string recognition_sequence_ ;  
vector<string> enzyme_acronyms_;
```

Other than the big-five (note that you can use the defaults for all of them), you have to add the following:

- A constructor `SequenceMap(const string &a_rec_seq, const string &an_enz_acro)`, that constructs a SequenceMap from two strings (note that now the vector `enzyme_acronyms_` will contain just one element, the `an_enz_acro`).
- `bool operator<(const SequenceMap &rhs) const`, that operates based on the regular string comparison between the `recognition_sequence_strings` (this will be a one line function).
- Overload the `operator<<` for SequenceMap.
- `void Merge(const SequenceMap &other_sequence)`. This function assumes that the object's `recognition_sequence_` and `other_sequence.recognition_sequence_` are equal to each other. The function `Merge()` merges the `other_sequence.enzyme_acronym_` with the object's `enzyme_acronym_`. The `other_sequence` object will not be affected.

This class (which is non-templated) will be used in the following programs. First test it with your own test functions to make sure that it operates correctly.

Part 2

Introduction to the problem

For this assignment you will receive as input two text files, `rebase210.txt` and `sequences.txt`. After the header, each line of the database file `rebase210.txt` contains the name of a restriction enzyme and possible DNA sites the enzyme may cut (cut location is indicated by a ') in the following format:

enzyme_acronym/recognition_sequence/.../recognition_sequence//

For instance the first few lines of `rebase210.txt` are:

```
AanI/TTA'TAA//  
AarI/CACCTGCNNNN'NNNN/'NNNNNNNNGCAGGTG//  
AasI/GACNNNN'NNGTC//  
AatII/GACGT'C//  
AbsI/CC'TCGAGG//  
AccI/GT'MKAC//  
AccII/CG'CG//  
AccIII/T'CCGGA//  
Acc16I/TGC'GCA//  
Acc36I/ACCTGCNNNN'NNNN/'NNNNNNNNGCAGGT//  
...
```

That means that each line contains one enzyme acronym associated with one or more recognition sequences. For example on line 2:

The enzyme acronym `AarI` corresponds to the two recognition sequences `CACCTGCNNNN'NNNN` and `'NNNNNNNNGCAGGTG`.

Part 2(a) (45 points)

You will create a parser to read in this database and construct a search tree (either a regular BST or an AVL tree). For each line of the database and for each recognition sequence in that line, you will create a new `SequenceMap` object that contains the recognition sequence as its `recognition_sequence_` and the enzyme acronym as the only string of its `enzyme_acronyms_`, and you will insert this object into the tree. This is explained with the following *pseudo code*:

```
Tree<SequenceMap> a_tree;
string db_line;
// Read the file line-by-line:
while (GetNextLineFromDatabaseFile(db_line)) {
    // Get the first part of the line:
    string an_enz_acro = GetEnzymeAcronym(db_line);
    string a_reco_seq;
    while (GetNextRegocnitionSequence(db_line, a_rego_seq){
        SequenceMap new_sequence_map(a_reco_seq, an_enz_acro);
        a_tree.insert(new_sequence_map);
    } // End second while.
} // End first while.
```

In the case that the `new_sequence_map.recognition_sequence_` equals the `recognition_sequence_` of a node X in the tree, then the search tree's `insert()` function will call the `X.Merge(new_sequence_map)` function of the existing element. This will have the effect of updating the `enzyme_acronym_` of X. Note, that this will be part of the functionality of the `insert()` function. The `Merge()` will only be called in case of duplicates as described above. Otherwise, no `Merge()` is required and the `new_sequence_map` will be inserted into the tree.

To implement the above, write a test program named **query_tree** which will use your parser to create a search tree and then allow the user to query it using a recognition sequence. If that sequence exists in the tree then this routine should print all the corresponding enzymes that correspond to that recognition sequence.

Your programs should run from the terminal as follows:

query_tree <database file name> <flag>

<flag> should be "BST" for binary search tree, and "AVL" for AVL tree.

For example you can write on the terminal:

`./query_tree rebase210.txt BST`

The user should enter THREE strings (supposed to be recognition sequences) for instance:

CC'TCGAGG
TTA'TAA
TC'C

Your program should print in the standard output their associated enzyme acronyms. In the above example the output will be

AbsI
AanI PsiI
Not Found

I will test it with a file containing three strings and run your code like that:

```
./query_trees rebase210.txt BST < input_part2a.txt
```

```
./query_trees rebase210.txt AVL < input_part2a.txt
```

Part2(b) (45 points)

Next, create a test routine named **test_tree** that does the following in the sequence described below:

1. Parses the database and construct a search tree (this is the same as in Part2(a)).
2. Prints the number of nodes in your tree n .
3. Computes the average depth of your search tree, i.e. the internal path length divided by n .
 - a. Prints the average depth.
 - b. Prints the ratio of the average depth to $\log_2 n$. E.g., if average depth is 6.9 and $\log_2 n = 5.0$, then you should print $\frac{6.9}{5.0} = 1.38$.
4. Searches (find()) the tree for each string in the sequences.txt file. Also counts the total number of recursive calls for all executions of find().
 - a. Prints the total number of successful queries (number of strings found).
 - b. Prints the average number of recursion calls, i.e. #total number of recursion calls / number of queries.
5. Removes every other string in sequences.txt from the tree. Also counts the total number of recursion calls for all executions of remove().
 - a. Prints the total number successful removes.
 - b. Prints the average number of recursion calls, i.e. #total number of recursion calls / number of remove calls.
6. Redo steps 2 and 3:
 - a. Prints number of nodes in your tree.
 - b. Prints the average depth.
 - c. Prints the ratio of the average depth to $\log_2 n$.

The output of Part2(b) should be of the exact form:

2: <integer>

3a: <float>

3b: <float>

4a: <integer>

4b: <float>

5a: <integer>

5b: <float>

6a: <integer>

6b: <float>

6c: <float>

If you didn't complete a step, just print after the step number: Not Done

For both Part2(a) and Part2(b) you **must** write the test routine using templates so each tree can be used interchangeably. Templated test routines are provided in the provided code. The trees should have identical interfaces.

Your program should run from the terminal as follows:

```
test_tree <database file name> <sequences file name> <flag>
```

<flag> should be "BST" for binary search tree, and "AVL" for AVL tree.

For example you can write on terminal

```
./test_tree rebase210.txt sequences.txt AVL
```

You will be given a mandatory Makefile, along with some code to start.
