

CS119 – Lab 8
Due Date: April 10

Purpose: Not all data can or should be represented linearly in a list. Hierarchical data abounds and is represented in the form of a tree. Examples from real life are family trees, the table of contents of a book, computer files system folders and subfolders, etc. We will look at a particular example involving data compression where a tree representation is useful.

Knowledge: This lab will help you become familiar with the following content knowledge:

- How to use and manipulate trees
- How trees can represent the encoding of text

Task: Follow the steps in this lab carefully to complete the assignments. Copy the lab8 folder and complete the following assignments by writing functions in `Huffman.hs`.

An illustration of the use of binary trees is in the problem of data compression. Ordinarily, each character is represented by an 8-bit code. We can reduce the total number of bits required to code a text by replacing this fixed-length code with a coding scheme based on the frequency of occurrence of characters in the text. Characters that appear most frequently should have short codes whereas characters that appear infrequently can have longer codes. For example in the word "text" we could encode $t \rightarrow 0$, $e \rightarrow 10$, and $x \rightarrow 11$. Then the encoding of "text" would be 010110.

We have to be careful, however, that we choose the code so that it can be uniquely decoded. If we had chosen $t \rightarrow 0$, $e \rightarrow 10$, and $x \rightarrow 1$, then both "text" and "tee" would have the same code! Not good. To prevent this from happening we must choose the codes so that no code is a proper prefix of any other.

To construct an optimal code satisfying the prefix property, we will use a technique called Huffman coding (named after David Huffman). Each character is stored as the leaf in a binary tree in such a way that more frequently used characters are of lesser depth in the tree than less frequently used ones. The code of character is a sequence of 0's and 1's describing the path in the tree to the character, where a 0 represents a left branch and a 1 a right branch. Consider the following tree structure and tree:

```
data HTree = Leaf Char | Branch HTree HTree deriving Show

Branch (Branch (Leaf 'x') (Leaf 'e')) (Leaf 't')
```

In this tree, the character 'x' is coded by 00, 'e' by 01, and 't' by 1.

To build a Huffman tree we start with a list of characters along with their frequencies. For example:

```
[('g',8),('r',9),('a',11),('t',13),('e',17)]
```

We convert this list of pairs into a list of trees and then repeatedly combine the trees with the lightest weights until just one tree remains. The weight of a single leaf will be the weight of the character at that leaf. The weight of a binary node is the sum of the weights of its two subtrees. We will need another tree data type for this weighted tree:

```
data WeightedTree = Tip Int Char | Node Int WeightedTree WeightedTree
                  deriving Show
```

After the weighted tree is constructed we can simply remove the weights and get our HTree.

The code for construction of the weighted tree is given to you. Look through the code and trace through it with the given frequencies. Test it out and see if you get the weighted tree that you expect.

Assignment 1:

We will make our HTree as follows:

```
makeHTree :: [(Char,Int)] -> HTree
makeHTree x = unweight (makeWeightedTree x)
```

Write the function `unweight` which takes a weighted tree and converts it to an HTree by stripping off the weights.

Criteria for Success: Test it out by creating the `huffTree` which uses the weights in the example above. Draw out the tree from the expression that is printed and verify that it is a tree that you would expect to get.

Assignment 2:

Now that we have our HTree it should be fairly straightforward to decode a message. Simply traverse the tree making left branches for 0's and right branches for 1's until you get a leaf. The given character is produced and if you have more bits, repeat the process again starting at the root for the next character.

Write the function `decode :: HTree -> [Bit] -> [Char]`.

Criteria for Success Test it out your function with the HTree from Assignment 1 and the bit string `[1,1,0,1,1,1,0,0,0,0,1]`. You should be able to figure out whether the text produced is correct or not.

Assignment 3:

Encoding is not as easy since the tree is good for finding a character associated with a bit string but poor at finding the bit string associated with a given character. So we will write a function `transform` which will transform the tree into a table where we can look up the code for a given character. This table and transform function will be

```
type CodeTable = [(Char,[Bit])]

transform :: HTree -> CodeTable
transform (Leaf x) = [(x,[])]
transform (Branch t1 t2) = hufmerge (transform t1) (transform t2)
```

Write the function `hufmerge :: CodeTable -> CodeTable -> CodeTable` which takes two code tables and merges them, adding a zero bit to the front of all the codes coming from the first table, and a one bit to the front of all the codes coming from the second table. Test it out by doing a transform of your HTree.

Criteria for Success: Use the `transform` function on the tree that you created previously. Take a look at the table that it is produced and verify that each letter has the correct code.

Assignment 4:

Write the function `codeLookup :: Char -> CodeTable -> [Bit]` which looks up the bit string for a given character in the CodeTable.

Criteria for Success: Perform `codeLookup 'e' (transform huffTree)` and verify that you get the correct bit string for that letter.

Assignment 5:

Write the function `encode :: HTree -> [Char] -> [Bit]` which uses a CodeTable to encode a string into a bits.

Criteria for Success: Encode the string "great" and verify that you get the bit string in Assignment 2.

Submit your `Huffman.hs` file in Canvas for grading.