

Homework 4

SNU 4190.210 Fall 2012

Chung-Kil Hur

due: 10/26(Sun) 24:00

Exercise 1 “Huffman Code”

Let’s try to devise a way to encode sentences, where a sentence is a finite sequence of words. For example, suppose we have four words: “I”, “me”, “you” and “know”. One can make many sentences using these words, such as “I know you” and “you know I know you know me”.

We will represent such sentences using 0 and 1. First, we assign to each word binary code consisting of 0 and 1. Then, a sentence is simply encoded as the sequential composition of the code for each word in the sentence. There are two ways to do this.

1. Fixed-length encoding: Since there are four words, each word can be encoded using two bits: 00 for “I”, 01 for “me”, 10 for “you” and 11 for “know”. Then, “you know I know you know me” is encoded as 10110011101101, which is 16 bit long. Decoding such code is easy: just decode each two bits one by one.
2. Variable-length encoding: If we know the frequency of each word, we can encode a sentence more efficiently by assigning short code for more frequent words.

Suppose that “know” is the most frequent and “you” is the second. Then we can encode “know” as 0, “you” as 10, “I” as 110 and “me” as 111. This way, we can encode “you know I know you know me” as 1001100100111, which is 13 bit long.

Decoding such code is also easy: just decode from left to right. For example, in 1001100100111, the only first possible code is 10, which is for

“you”, and the only second possible code is 0, which is for “know”, and so on. We can easily see that 1001100100111 is uniquely decoded into “you know I know you know me”.

Why is the decoding easy? The reason is because there is no word code that is a prefix of another word code. The word codes 0, 10, 1110 and 111 satisfy this property, which is called *prefix-free*. Such a variable-length encoding is used in data compression software such as JPEG, MPEG and ZIP.

This assignment is to write a function `vlenencode` that finds the optimal prefix-free variable-length code. Given a list of pairs of a word and its frequency, `vlenencode` should return a list of pairs of a word and its code. A word is a string, a frequency is an integer and code is a list of 0 and 1.

In 1951, David Huffman, as a graduate student at MIT, found an optimal solution for the variable-length encoding problem as a term project in the course “Information Theory”. With this result, he outperformed Prof. Robert Fano, who taught the course and had worked on the same problem with Claude Shannon, the inventor of information theory.

You can also do it as Huffman did. The following hints might be enough. Try to solve it by yourself before googling the solution.

- Use a binary tree to generate prefix-free code. Words are assigned to leaves of a binary tree, and each branch in the tree determines 0 or 1.

Define and use the following functions to implement binary trees.

```
leaf : int × string → tree
node : int × tree × tree → tree
isleaf? : tree → bool
leafval : tree → int
leafstr : tree → string
nodeval : tree → int
leftsub : tree → tree
rightsub : tree → tree
```

- Build a binary tree, taking into account the frequency of each word. For this, store the total frequency of words in a sub-tree in the root of the sub-tree.

□

Exercise 2 “SKI combinator calculus“

SKI expressions E are inductively defined as follows:

$$\begin{array}{l} E \rightarrow S \mid K \mid I \\ \quad \mid x \quad \quad \text{variables} \\ \quad \mid (E E) \end{array}$$

Examples are

$$K, (I x), (S ((K x) y)), (((S K) K) x).$$

The evaluation rules of SKI calculus is as follows.

$$\begin{array}{l} (I E) \rightarrow E \\ ((K E) E') \rightarrow E' \\ (((S E) E') E'') \rightarrow ((E E'') (E' E'')) \end{array}$$

For example, $((S K) I) x$ evaluates as follows:

$$(((S K) I) x) \rightarrow ((K x) (I x)) \rightarrow x$$

Define a function `execute` that takes an SKI expression and prints out a sequence of SKI expressions that occur during the evaluation.

$$\text{execute} : E \rightarrow \text{void}$$

Note that there are several possible evaluations of the same expression. For example, $((K x) (I y))$ can evaluate as follows:

$$((K x) (I x)) \rightarrow x$$

Or,

$$((K x) (I x)) \rightarrow ((K x) x) \rightarrow x.$$

In such a case, `react` is allowed to follow just one of the possible executions.

Implement functions with the following interface as a library for SKI cal-

culus, and use the library to define `react`.

```
S : E
K : E
I : E
v : string → E      (* construct an SKI expression consisting of a variable *)
a : E × E → E      (* construct an SKI expression of the form (E1 E2) *)
isS? : E → bool
isK? : E → bool
isI? : E → bool
isv? : E → bool
isa? : E → bool
var : E → string    (* When the argument is a variable, return its name *)
al : E → E          (* When the argument is (E1 E2), return E1 *)
ar : E → E          (* When the argument is (E1 E2), return E2 *)
pprint : E → void   (* pretty-print the given SKI expression *)
```

where `pprint` should print SKI expressions in the way that TA will tell you. \square