

Arbres préfixes

OPTION INFORMATIQUE - TP n° 1.5 - Olivier Reynet

On se dote d'un type pour représenter les arbres préfixe pour élaborer un code de Huffman :

```
type btree = Leaf of char
           | Node of int * btree * btree;;
```

Ces arbres sont constitués de feuilles étiquetées par des types char (les symboles à encoder) et des nœuds étiquetés par des types int (les nombre d'occurrences).

Pour transformer une chaîne de caractères en liste de caractères, on pourra utiliser le code suivant :

```
let string_to_char_list s = List.of_seq (String.to_seq s)
```

1. Écrire une fonction de signature `occurences : string -> (btree * int)list` dont le paramètre est un message à encoder. Cette fonction renvoie la liste des couples feuilles et nombre d'occurrences associés à chaque symbole. On utilisera une table de hachage pour effectuer le décompte des occurrences de chaque caractère dans le message. Cette table sera transformée en liste par la commande `Hashtbl.fold (fun k v acc -> (Leaf k, v):: acc)dict []`.

Par exemple, `occurences "Hello"` renvoie

```
[(Leaf 'l', 2); (Leaf 'o', 1); (Leaf 'e', 1); (Leaf 'H', 1)].
```

Pour construire l'arbre d'Huffman, il faut Implémenter une file de priorités. On se propose de réaliser celle-ci en utilisant :

- une liste dont les éléments sont des couples de type `(btree * int)`,
 - une fonction `insert_elem` qui insère un élément dans une telle liste, au bon endroit, c'est à dire d'après l'entier qui code le nombre d'occurrences du symbole : les plus petits nombres en tête.
 - une fonction `insert_sort` qui implémente le tri par insertion pour une liste de type `(btree * int)list`.
2. Écrire une fonction de signature `compare : 'a * int -> 'b * int -> int` dont les paramètres sont deux couples (symbole, occurrences) de la liste et qui renvoie le résultat de `n1-n2` si `n1` et `n2` sont les occurrences associés aux symboles des deux couples.
 3. Écrire une fonction récursive de signature `insert_elem : ('a * int)list -> 'a * int -> ('a * int)list` qui insère un couple (symbole, occurrences) à la bonne place dans la liste passée en paramètre.
 4. Écrire une fonction récursive de signature `insert_sort : ('a * int)list -> ('a * int)list` qui implémente le tri par insertion.
 5. Écrire une fonction de signature `merge : btree * int -> btree * int -> btree * int` qui fusionne les deux sous-arbres en un arbre `btree` et renvoie le couple constitué du nouveau nœud et de la somme des occurrences de sous-arbres. Par exemple, `merge (Leaf 'a', 5) (Leaf 'b', 9)` renvoie `(Node (14, Leaf 'a', Leaf 'b'), 14)`.
 6. Écrire une fonction de signature `huffmann_tree : (btree * int)list -> btree` qui construit l'arbre de Huffman associé à une liste de couples (symbole, occurrences). Par exemple :

```

let q = [(Leaf 'a',5); (Leaf 'b',9); (Leaf 'c',12); (Leaf 'd',13); (Leaf 'e',16);
         (Leaf 'f',45)];;
let ht = huffmann_tree q;;
(* Node (100, Leaf 'f',
        Node (55, Node (25, Leaf 'c', Leaf 'd'),
                Node (30, Node (14, Leaf 'a', Leaf 'b'), Leaf 'e')))) *)

```

7. Écrire une fonction de signature `h_decode : btree -> string -> string` qui décode un message donné sous la forme d'une chaîne de caractères ne comportant que des 0 et des 1. Cette fonction renvoie la chaîne de caractères correspondant au message initial. On pourra utiliser l'arbre de Huffman défini à la question précédente et décoder "0100111".
8. Écrire une fonction de signature `encode_map : btree -> (char, string)Hashtbl.t` dont le paramètre est un arbre de Huffman et qui renvoie une table de hachage associant chaque symbole à son encodage binaire sous la forme d'une chaîne de caractères. On pourra s'appuyer sur les fonction `String.to_seq` et `List.of_seq`.
9. Écrire une fonction de signature `h_encode : (char, string)Hashtbl.t -> string -> string` dont les paramètres sont une table d'encodage et un message sous la forme d'une chaîne de caractères et qui renvoie la chaîne binaire qui encode le message. On pourra s'appuyer sur les fonction `String.to_seq` et `List.of_seq`.
10. Écrire une fonction de signature `compression_rate : string -> float` qui calcule le taux de compression d'une chaîne de caractère encodée par l'arbre préfixe de Huffman. Le taux de compression se calcule : $t = 1 - \frac{c}{8 \times n}$, si c est la longueur de la chaîne binaire compressée et n la longueur de la chaîne du message de départ

(R) La limite théorique de compression est atteignable avec un code d'Huffman. On peut calculer le nombre de bits moyen nécessaire pour encoder un symbole. Celui-ci s'écrit :

$$H = - \sum_{i=0}^{n-1} p_i \log_2 p_i \quad (1)$$

où n est le nombre de symboles et p_i la probabilité d'apparition du symbole i .