

COMP 212 Spring 2023

Homework 05

This homework will focus on lists, trees, and work-span analysis. You must write purposes and tests for all functions on this assignment.

1 Trees

The type `tree`, with constructors `Empty` and `Node(l,x,r)` represents binary trees of integers with data stored only in the internal nodes. Here are some definitions about trees:

- The tree `Empty` has *depth* 0. The tree `Node(l,x,r)` has *depth* d if and only if
 1. l has depth d_l
 2. r has depth d_r
 3. $d = \max(d_l, d_r) + 1$.
- The tree `Empty` has *size* 0. The tree `Node(l,x,r)` has *size* s if and only if
 1. l has size s_l
 2. r has size s_r
 3. $s = s_l + s_r + 1$
- The tree `Empty` is *balanced*. The tree `Node(l,x,r)` is *balanced* if and only if
 1. l is balanced
 2. r is balanced
 3. l has depth d_l , r has depth d_r and $|d_l - d_r| \leq 1$.
- The tree `Empty` is *sorted*. The tree `Node(l,x,r)` is *sorted* if and only if
 1. l is sorted
 2. r is sorted
 3. For every node `Node(l1,x1,r1)` in l , $x_l < x$
 4. For every node `Node(lr,xr,rr)` in r , $x_r \geq x$

An expression $e : \text{tree}$ is sorted iff it is equivalent to a value that is sorted.

These definitions are implemented in `hw05-lib.sml`, with a few other helper functions. Download this file and put it in the same folder as your `hw05.sml` file. You should feel free to write your tests in terms of these functions.

- `depth : tree -> int` computes the depth of its argument.
- `size : tree -> int` computes the size of its argument.
- `isbalanced : tree -> bool` evaluates to true if and only if its argument is balanced.
- `issorted : tree -> bool` evaluates to true if and only if its argument is sorted.
- `tolist : tree -> int list` computes a flattening of its argument into a list, as given in lab.
- `fromlist : int list -> tree` computes a balanced tree from a list—very useful for testing, but do **not** use it in any of your solutions, or they will not have the right span.
- `treeeq : tree * tree -> bool` tests whether two trees are equal

2 Contains

Task 2.1 (7 pts). Write a function `contains : tree * int -> bool` such that for any sorted tree `t`, `contains(t,i)` returns `true` if `i` is an element of `t` and returns `false` if not. For full credit, your solution should have $O(\log n)$ work and span when t is a balanced tree of size n .

The `listToTree` function from lab is called `fromlist : int list -> tree` in the homework code, and can be used to write tests.

3 Tree Induction

Recall the `size` and `depth` functions on trees:

```
(* size t computes the natural number which
   represents how many Nodes are in a tree *)
fun size (t : tree) : int =
  case t of
    Empty => 0
  | Node (l, x, r) => 1 + size l + size r

(* depth t computes the natural number which
   represents how many levels are in a tree *)
fun depth (t : tree) : int =
  case t of
    Empty => 0
  | Node (l, x , r) => 1 + max (depth l, depth r)
```

Task 3.1 (15 pts). Prove the following relationship:

Theorem 1. *For all trees t , $\text{depth } t \leq \text{size } t$.*

The function `max(x,y)` returns the maximum of x and y ; i.e. if x is bigger it returns x and if y is bigger it returns y . You may want to use some of the following properties of `max`:

- $x \leq \text{max}(x, y)$ and $y \leq \text{max}(x, y)$
- $\text{max}(x, y) \leq z$ if $x \leq z$ and $y \leq z$
- $\text{max}(x, y) \leq \text{max}(x', y')$ if $x \leq x'$ and $y \leq y'$

4 QuickSort

In last week's homework, you implemented `QUICKSORT` on lists. As we've discussed in lecture, there is not a lot to be gained by using parallel sorting algorithms on lists: there are dependencies inherent in the structure of a list that get in the way of real parallelism.

In that spirit, you'll now implement `QUICKSORT` on trees. Assuming the pivots yield subproblems of equal size (which can be achieved using randomness), this algorithm will have $O(n \log n)$ work and $O((\log n)^3)$ span. The logarithmic span means significant speedups can be gained by running the algorithm in parallel.

We'll represent trees with the `tree` type defined at the beginning of this assignment. In specs, we will say that `x` is an element of a tree `t` when `Node(...,x,...)` appears somewhere in `t`.

4.1 Combine

First, we will need a function to combine two trees into one. Unlike `merge` from lecture, we will not require that the inputs are sorted, but we will also not ensure that the outputs are sorted, or that the outputs are in the same order as in the original trees. Because this function will be used prior to sorting, the elements can be in any order we choose. This means the following code suffices:

```
fun combine (t1 : tree, t2 : tree) : tree =  
  case t1 of  
    Empty => t2  
  | Node(l1,x1,r1) => Node(combine(l1,r1),x1,t2)
```

You may wish to draw `combine`'s output on some examples to understand how it works. More formally, this code meets the following specification:

- **Functionality:** For all trees `T1` and `T2`, `combine (T1,T2)` is valuable, and contains every element of `T1` and every element of `T2` and no other elements.
- **Depth:** For the analysis of `quicksort`, we need the following bound on the depth of `combine`'s result:

Lemma 1. *For all values `t1 t2:tree`,
 $\text{depth } (\text{combine}(t1,t2)) \leq 1 + \max(\text{depth } t1, \text{depth } t2)$.*

- **Running-time:** Let d_1 be the depth of `T1`, d_2 be the depth of `T2`. The work and span of `(combine (T1,T2))` are both $O(d_1)$.

4.2 Filter

Task 4.1 (15 pts). Implement a tree analogue of `filter_less` and `filter_greatereq`:

```
filter_less : tree * int -> tree
filter_greatereq : tree * int -> tree
```

Your implementation must satisfy the following specs:

- **Functionality:** If T is a value of type `tree`, p is a value of type `int`, then:
 - `filter_less(T,p)` contains all and only the elements of T that are less than p .
 - `filter_greatereq(T,p)` contains all and only the elements of T that are greater than or equal to p .
- **Depth:** For all $T:\text{tree}, p:\text{int}$, $\text{depth}(\text{filter}(T,p)) \leq \text{depth } T$.
- **Running-time:** If d is the depth of a tree T , your implementation of each `(filter(T,p))` should have $O(d^2)$ span. On a balanced tree, your implementation of each `filter` should have $O(n)$ work, where n is the size of the tree.¹

4.3 Quicksort

Task 4.2 (15 pts). Finally, put all the pieces together to write

```
quicksort_t: tree -> tree
```

which implements QUICKSORT values of type `tree`.

- **Functionality:** `quicksort_t T` is sorted and contains all and only the elements of T .
- **Running-time:** If T is a balanced tree with size n , `(quicksort_t T)` should have $O(n \log n)$ work and $O((\log n)^3)$ span, assuming the pivots yield balanced subproblems.

You may want to use the `fromlist : int list -> tree` and `issorted` functions to test your implementation of `quicksort_t`.

5 Balancing

To `mergesort` trees, we needed to *rebalance* a tree after manipulating it. Rebalancing takes a tree that is not necessarily balanced, and computes a balanced tree with the same elements.

We have provided most of an implementation of a simple rebalancing algorithm in the handout. The key helper function is unimplemented. You will implement this helper and then analyze the complexity of `rebalance`.

In all of the tasks, you should assume that the function `size : tree -> int`, which computes the size of a tree, runs in constant time on all inputs. This happens to be obviously false. However, it's easy to make binary trees whose size can be computed in constant time by storing the size at each node—so this is a relatively harmless lie.

Task 5.1 (15 pts). Implement the function

¹If you use the `tree` method to try to prove this, you will run into a sum that we have not yet seen in the course; see the next problem for its big-O.

```
takeanddrop : tree * int -> tree * tree
```

`takeanddrop(T,i)` separates a tree `T` into “left” and “right” subtrees, `T1` and `T2` respectively. `T1` contains the leftmost `i` elements of `T`, in their original order, and `T2` the remaining elements, also in their original order. For example, if we define

```
val test =
  Node
    (Node (Node (Empty,1,Empty),
              2,
              Node (Empty,3,Empty))),
    4,
    Node (Node (Node (Empty,5,Empty),
                  6,
                  Empty)))
```

then we have

```
takeanddrop (test,3) ==
  (Node (Node (Node (Empty,1,Empty),2,Node (Empty,3,Empty)),
    Node (Empty,4,Node (Node (Empty,5,Empty),6,Empty)))
```

More formally, suppose `T` is any tree, and $0 \leq i \leq \text{size } T$. Then `takeanddrop (T,i)` evaluates to a pair of trees `(T1,T2)` such that

- $\max(\text{depth } T1, \text{depth } T2) \leq \text{depth } T$
- $\text{size } T1 \cong i$
- $(\text{tolist } T1) @ (\text{tolist } T2) \cong (\text{tolist } T)$

This last condition ensures that `T1` contains the leftmost elements, and that the elements of `T1` and `T2` are in the appropriate order.

If d is the depth of `T` then your implementation of `(takeanddrop (T,i))` must have $O(d)$ work and span.

Hint: use the `splitAt` function from mergesorting trees as a model; the difference is that instead of splitting based on the values stored in the tree, here we are splitting based on the number of elements in the tree.

Task 5.2 (18 pts). Your implementation of `takeanddrop` is necessary for the helper function `halves`, which is used by `rebalance`; see the starter code. The following tasks ask you to analyze these functions:

1. Give a recurrence that describes the work of your implementation of `takeanddrop`, $W_{\text{takeanddrop}}(d)$, in terms of the **depth** d of the input tree. Argue that $W_{\text{takeanddrop}}(d)$ is $O(d)$.
2. Give a recurrence that describes the span of your implementation `takeanddrop`, $S_{\text{takeanddrop}}(d)$, in terms of the **depth** d of the input tree. Argue that $S_{\text{takeanddrop}}(d)$ is $O(d)$

Note: the remaining tasks will be graded assuming that $W_{\text{takeanddrop}}(d)$ and $S_{\text{takeanddrop}}(d)$ are $O(d)$, even if that is not true for your code, or if your recurrence above is incorrect.

3. Give a recurrence that describes the work of **halves**, $W_{\text{halves}}(d)$, in terms of the **depth** of the input tree. Give a tight big-O bound for $W_{\text{halves}}(d)$.
4. Give a recurrence that describes the span of **halves**, $S_{\text{halves}}(d)$, in terms of the **depth** of the input tree. Give a tight big-O bound for $S_{\text{halves}}(d)$.
5. Give a recurrence that describes the work of **rebalance**, $W_{\text{rebalance}}(n)$, in terms of the **size** n of the input tree. You should assume that the input is roughly balanced—that is to say, there exists some constant c such that the depth of the input is $c \log n$. This will be true when **rebalance** is called from **mergesort**, because the merging will only have unbalanced the tree by a known amount.

Give a tight big-O bound for $W_{\text{rebalance}}(n)$. Show your work using a closed form and/or sum.

6. Give a recurrence that describes the span of **rebalance**, $S_{\text{rebalance}}(n)$, in terms of the **size** of the input tree. You should assume that the input is roughly balanced—that is to say, there exists some constant c such that the depth of the input is $c \log n$. This will be true when **rebalance** is called from **mergesort**, because the merging will only have unbalanced the tree by a known amount.

Give a tight big-O bound for $S_{\text{rebalance}}(n)$. Show your work using a closed form and/or sum.

The recurrences for the later tasks should be defined in terms of the recurrences defined in the earlier tasks for the helper functions.

You may use the following tight bounds:

$$\begin{array}{ll} \log n + \log \frac{n}{2} + \log \frac{n}{4} + \log \frac{n}{8} + \dots + 1 & \text{is } O(\log n)^2 \\ \log n + 2 \log \frac{n}{2} + 4 \log \frac{n}{4} + 8 \log \frac{n}{8} + \dots \text{ (with } \log n \text{ terms)} & \text{is } O(n) \end{array}$$

6 NON-COLLABORATIVE CHALLENGE PROBLEM: Burrows-Wheeler Transform

Remember that non-collaborative challenge problems are to be done independently. You are not allowed to communicate with anyone about the problems, except to ask the instructor or TAs clarification questions (not hints). Additionally, you are not allowed to search for help on the specific problem from any sources besides the course materials. You are free to ask clarification questions about the concepts involved — in this case, recurrences and big- O .

In general, *non-lossy data compression* is about making a smaller representation of some data without losing any information. E.g. this is what zip files do — they make a smaller file to download that can be “decompressed” into the original data.

In this problem, you will analyze the running time of an implementation of a particular kind of data compression for text documents, called the Burrows-Wheeler Transform followed by run-length encoding.

The code for this algorithm is in a separate file linked from the assignments page (use the password for the lab solutions).

We'll say that a `(char list) list` is *rectangular* iff it is a list of n lists, each of which has the same length m . This means that if you write it out with one element on each line it will look like a rectangle, for example

```
[ ["a","b","c"],  
  ["b","c","a"],  
  ["c","a","b"] ]
```

is rectangular but

```
[ ["a","b"],  
  ["b"],  
  ["c","a","b"] ]
```

is not.

Task 6.1 (22 pts). Give a tight big-O bound for the work of each of the following functions. For each function, clearly state what input sizes you are analyzing the function in terms of. Note that some of the later functions are easier to analyze than the earlier ones, so you should try all of them even if you get stuck on some. You can show recurrences/closed forms for partial credit in case your big-O is incorrect, but if your big-O is correct you will get full credit.

- `rotations_helper`. You can assume that `fastReverse` and `append (@)` have the work we analyzed in class.
- `rotations`
- `list_lt`. You can assume that the two input lists have the same length.
- `merge`. You can assume the input is rectangular.
- `mergesort`. You can assume `split` has the running time we analyzed in class, and that the input list is rectangular.
- `last`
- `all_last`s. You can assume that the input list is rectangular.
- BWT
- `lasHelp`. You will receive full credit if your bound is not as tight as possible, as long as this doesn't change the overall big-O for `compress`.

- `run_length_encode`. You can assume that string-append \sim has the same work as the corresponding list append. You will receive full credit if your bound is not as tight as possible, as long as this doesn't change the overall big-O for `compress`.
- `compress`

Because this problem is a little open-ended, there are no boxes for it in `hw05-written.pdf`. You can either append it to the end of your written solution or upload a separate `hw05-challenge.pdf` to your Google Drive.