Chapter 6

Greedy Algorithms

6.1 Storing the Huffman codeword tree T on a file

We briefly describe a possible way of storing a compact representation of the tree T corresponding to the optimal prefix code on a file. Such a representation is essential and must be stored together with the encoding of a given text file so that a decompression routine can reconstruct it and use it for decoding.

Define the *preorder sequence* pre(T) of a given full binary tree T as a string containing all nodes of T ordered as follows:

- 1. If T is a single-node tree, then pre(T) = (root(T)), where root(T) denotes the only node of T (which is also its root).
- 2. If T contains more than one node, let left(T) and right(T) denote, respectively, its left and right subtree (note that these subtrees are both nonempty since T is full). Then $pre(T) = \langle root(T), pre(left(T)), pre(right(T)) \rangle$

We will represent the codeword tree T in the file by means of its preorder sequence. Recall that internal nodes do not have to carry any information, while for each leaf f we must store the character char(f) associated with it. Therefore, in the preorder sequence, we identify each internal node with a single bit, 0, while a leaf f is identified with the pair of symbols 1, char(f). Under this representation, the preorder sequence associated to the optimal tree T of figure 16.4.(b) (CLRS, page 387) is

 $\langle 0, 1, ASCII(a), 0, 0, 1, ASCII(c), 1, ASCII(b), 0, 0, 1, ASCII(f), 1, ASCII(e), 1, ASCII(d) \rangle$

Note that, since the tree t is full, in the sequence there are always n-1 0's and n pairs

 $1, \operatorname{char}(f).$

When decompressing, it is very simple to reconstruct T from pre(T). Indeed, it is sufficient to create the root node with the statement $r \leftarrow new_node()$ and then call the following routine with parameter r:

```
\begin{array}{l} \mbox{BUILD_TREE}(r) \\ b \leftarrow \mbox{read}() \\ \mbox{if } (b=0) \\ \mbox{then} \\ s \leftarrow \mbox{new\_node}() \\ \mbox{left}[r] \leftarrow s \\ \mbox{BUILD\_TREE}(\mbox{left}[r]) \\ t \leftarrow \mbox{new\_node}() \\ \mbox{right}[r] \leftarrow t \\ \mbox{BUILD\_TREE}(\mbox{right}[r]) \\ \mbox{else} \left\{ \begin{array}{c} b=1 \end{array} \right\} \\ c \leftarrow \mbox{read}() \\ \mbox{char}[r] \leftarrow \mbox{ASCII}(c) \\ \mbox{left}[r] \leftarrow \mbox{nil} \\ \mbox{right}[r] \leftarrow \mbox{nil} \\ \mbox{return} \end{array} \right.
```

Exercise 6.1 Let $S = \{[a_1, b_1], [a_2, b_2], \dots, [a_n, b_n]\}$ be a set of closed intervals on the real line. We say that $C \subseteq S$ is a *covering subset* for S if, for any interval $[a, b] \in S$, there exists an interval $[a', b'] \in C$ such that $[a, b] \subseteq [a', b']$ (that is, $a \ge a'$ and $b \le b'$).

- (a) Write a greedy $O(n \log n)$ algorithm that, on input S, returns a covering subset C^* of minimum size.
- (b) Prove the greedy choice property for the above algorithm.
- (c) Prove the optimal substructure property for the problem.

Answer:

(a) We first sort the intervals in nondecreasing order of their left endpoint. Ties are broken as follows: if $a_i = a_j$ then $[a_i, b_i]$ precedes $[a_j, b_j]$ in the sorted sequence if $b_i > b_j$. Otherwise, if $b_i < b_j$, $[a_j, b_j]$ precedes $[a_i, b_i]$ (note that it cannot be the case that $b_i = b_j$ or the two intervals would be the same). In other words, ties are broken in favor of the larger interval. Once the intervals are so sorted, we perform a linear scan of the intervals. The greedy choice is to select the first interval. At the *i*-th iteration, let $[a_k, b_k]$ denote the last interval that was selected. Interval $[a_i, b_i]$ is discarded if it is contained in $[a_k, b_k]$, and is selected otherwise.

We assume that the intervals are given in input stored into an array of records S, with S[i].left = a_i and S[i].right = b_i . In the code, procedure SORT(S) performs the required sorting of the intervals, so that, on its termination, S[i] stores the *i*-th interval of the sorted sequence. Note that a trivial adaptation of MERGE_SORT suffices to yield an $O(n \log n)$ running time for SORT. The algorithm follows:

$$\begin{split} \text{MIN-COVERING}(S) \\ n &\leftarrow \text{length}(S) \\ \text{SORT}(S) \\ C[1] &\leftarrow S[1] \\ j &\leftarrow 2 \\ \{ \text{ j stores the index of the first empty location in vector C } \} \\ k &\leftarrow 1 \\ \{ \text{ k stores the index of the last selected interval } \} \\ \text{ for } i &\leftarrow 2 \text{ to n do} \\ \text{ if $S[i].right > S[k].right } \\ \text{ then $C[j] \leftarrow S[i] \\ j &\leftarrow j + 1 \\ k &\leftarrow i \\ \{ S[i] \text{ is now the last selected interval } \} \\ \text{ return C} \end{split}$$

The running time $T_{M_{-C}}(n)$ of MIN₋COVERING(S) is clearly dominated by the time required by the initial sort, hence $T_{M_{-C}}(n) = O(n \log n)$.

(b) MIN_COVERING clearly exhibits the greedy choice property. In fact, *every* covering subset must contain S[1], since by construction S[1] is not contained in any other interval in S.

(c) Let $S = \{[a_1, b_1], [a_2, b_2], \ldots, [a_n, b_n]\}$ and assume that the intervals in S have been renumbered so to comply with the ordering enforced by procedure SORT above. Let $C^* = \{[a_1, b_1]\} \cup \overline{C}^*$ be a covering subset of minimum size for S. The optimal substructure property that the algorithm MIN_COVERING exploits is the following: letting k be the maximum index such that $b_j \leq b_1$, for $1 \leq j \leq k$ (note that $a_j \geq a_1$ also holds, because of the ordering), then \overline{C}^* must be an optimal solution to $\overline{S} = S - \{[a_j, b_j] : 1 \leq j \leq k\}$ (observe that \overline{S} could be empty). To prove this, let us first prove that \overline{C}^* is indeed a covering subset for \overline{S} . If $\overline{S} = \emptyset$, then \overline{C}^* is clearly empty, since $\{[a_1, b_1]\}$ is a covering subset for S. If \overline{S} is not empty, then \overline{C}^* must be a covering subset for all intervals in S not contained in $[a_1, b_1]$, since $C^* = \{[a_1, b_1]\} \cup \overline{C}^*$ is a covering subset. Therefore, \overline{C}^* must contain $[a_{k+1}, b_{k+1}]$, which is not covered by any other interval in S. Note, however, that \overline{S} might also contain intervals which are contained in $[a_1, b_1]$. Let $[a_s, b_s]$ be one of such intervals. First note that it must be s > k + 1, since the interval is in \overline{S} . Then, the following chain of inequalities is easily established:

$$a_1 \le a_{k+1} \le a_s \le b_s \le b_1 < b_{k+1},$$

hence $[a_s, b_s]$ is covered by $[a_{k+1}, b_{k+1}] \in \overline{C}^*$. To finish the proof it suffices to observe that if \overline{C}^* were not a covering subset of \overline{S} of minimum size, then C^* would not be a covering subset of minimum size for S.

Exercise 6.2 Let $C = \{1, 2, ..., n\}$ denote a set of n rectangular frames. Frame i has base d[i].b and height d[i].h. We say that Frame i encapsulates Frame j if $d[i].b \ge d[j].b$ and $d[i].h \ge d[j].h$. (Note that a frame encapsulates itself). An encapsulating subset $C' \subseteq C$ is such that for each $j \in C$ there exists $i \in C'$ such that i encapsulates j.

- (a) Design and analyze a greedy algorithm that, on input the vector d of sizes of a set C of n frames, returns a *minimum size* encapsulating subset $C' \subseteq C$ in $O(n \log n)$ time.
- (b) Prove the greedy choice property.
- (c) Prove the optimal substructure property.

Answer: We first sort the entries of vector d in nonincreasing order of their $d[\cdot].b$ field. Ties are broken as follows: if d[i].b = d[j].b then d[i] precedes d[j] in the sorted sequence if $d[i].h \ge d[j].h$. In other words, ties are broken in favor of the higher frame. Once the frames are so sorted, we perform a linear scan starting from the first frame in the sorted order. The greedy choice is to select the first frame. At the beginning of the *i*-th iteration, $i \ge 2$, let j < i denote the last frame that was selected. Frame *i* is discarded if it is encapsulated by frame *j*, and is selected otherwise.

In the code below, procedure SORT(d) performs the required sorting of the frames, so that, on its termination, d[i] stores the *i*-th frame of the sorted sequence. Note that a trivial adaptation of MERGE_SORT suffices to yield an $O(n \log n)$ running time for SORT. The algorithm follows:
$$\begin{split} &\text{MIN}_\text{ENCAPSULATING}_\text{SUBSET}(d) \\ &n \leftarrow \text{length}(d) \\ &\text{SORT}(d) \\ &C'[1] \leftarrow d[1] \\ &k \leftarrow 2 \\ \{ \ k \ \text{stores the index of the first empty location in vector } C' \ \} \\ &j \leftarrow 1 \\ \{ \ j \ \text{stores the index of the last selected frame } \} \\ &\text{for } i \leftarrow 2 \ \text{to } n \ \text{do} \\ & \text{if } d[i].h > d[j].h \\ & \text{then } C'[k] \leftarrow d[i] \\ &k \leftarrow k + 1 \\ &j \leftarrow i \\ &\{ \ i \ \text{is now the last selected frame } \} \\ &\text{return } C' \end{split}$$

The running time $T_{M_E_S}(n)$ of MIN_ENCAPSULATING_SUBSET(d) is clearly dominated by the time required by the initial sort, hence $T_{M_E_S}(n) = O(n \log n)$.

In the following, we assume that the frames in C have been renumbered so to comply with the ordering enforced by procedure SORT.

(b) Greedy Choice Consider an arbitrary optimal solution $C^* \subseteq C$. Such subset must contain a frame t which encapsulates Frame 1. Then, either t = 1, in which case C^* contains the greedy choice, or t > 1 and (1) $d[t].b \ge d[1].b$ and (2) $d[t].h \ge d[1].h$. In the latter case, due to the sorting, both (1) and (2) must be equalities, hence all frames encapsulated by Frame t are also encapsulated by Frame 1. Therefore $(C^* - \{t\}) \cup \{1\}$ is an optimal solution containing the greedy choice.

(c) Optimal Substructure Let $C^* = \{1\} \cup \overline{C}^*$ be an encapsulating subset of minimum size for C containing the greedy choice. The optimal substructure property featured by the problem is the following: let $j \leq n+1$ be the maximum value such that $d[i].h \leq d[1].h$ for $1 \leq i < j$. Then \overline{C}^* must be an optimal solution to $\overline{C} = C - \{i : 1 \leq i < j\}$ (observe that \overline{C} could be empty). To prove this, let us first prove that \overline{C}^* is indeed an encapsulating subset for \overline{C} . If $\overline{C} = \emptyset$, then \overline{C}^* is clearly an encapsulating subset. If \overline{C} is not empty, then \overline{C}^* must be an encapsulating subset for all frames in C which are not encapsulated by Frame 1, since $C^* = \{1\} \cup \overline{C}^*$ is an encapsulating subset for C. Therefore, \overline{C}^* must contain Frame j, which is not encapsulated by any other frame in C (indeed, \overline{C}^* could contain a frame j' with the same base and height as j, but then we can substitute j' with j in \overline{C}^* and proceed with the argument). Note, however, that \overline{C} might also contain frames which are encapsulated by Frame 1. Let s be one of such frames. First note that it must be $s \ge j$, since the frame is in \overline{C} . Then, $d[j].b \ge d[s].b$ (by the ordering) and

$$d[j].h > d[1].h \ge d[s].h$$

hence Frame s is also encapsulated by $j \in \overline{C}^*$.

To finish the proof it suffices to observe that if \overline{C}^* were not an encapsulating subset of \overline{C} of minimum size, then C^* would not be an encapsulating subset of C of minimum size. \Box

Exercise 6.3 Consider a variant of the Activity Selection Problem, where the input set of intervals $S = \{[s_1, f_1), ..., [s_n, f_n)\}$ is sorted by non decreasing values of the s_i 's, that is, $s_1 \leq s_2 \ldots \leq s_n$. As in the original problem, we want to determine a maximum set of pairwise disjoint activities.

- (a) Design an O(n) algorithm for the above problem.
- (b) Prove that the greedy choice and the optimal substructure properties hold.

Exercise 6.4 Let $X = \{x_1, x_2, \dots, x_N\}$ be a set of points on the real line.

- (a) Design and analyze a greedy algorithm to determine a minimum cardinality set I of closed intervals of unit length (that is, $i = [a, b] \in I \Rightarrow (b a) = 1$) such that, for any $x \in X$, there exists an interval $i \in I$ which contains x.
- (b) Prove that your algorithm has the greedy choice property.
- (c) Prove that the problem enjoys the optimal substructure property.