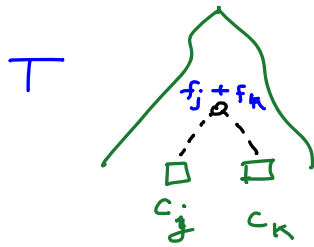


Huffman Codes - Continued

Claim: Let $C = \{(c_i, f_i) \mid 1 \leq i \leq n\}$ be an instance of the (Huffman coding) problem. Let f_j, f_k be two minimum frequencies ($j \neq k$) in C . Then there is an optimal solution (to C) with choice f_j, f_k .

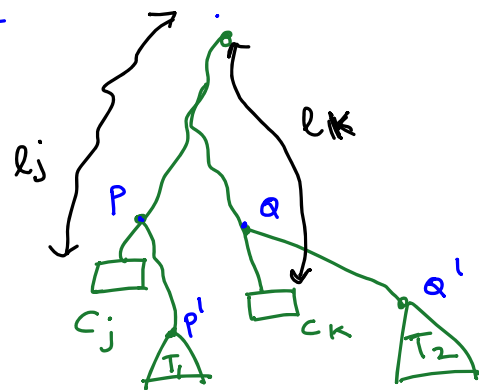
Proof: Let T be tree corresponding to codes of some optimal solution of C .



If leaves c_j, c_k are children of some internal node u ,

then this optimal solution can be seen as greedy choice + solution to the subproblem.
So we are done in this case.

T



l_j : length of the path from the root to c_j (length of the code word for c_j)

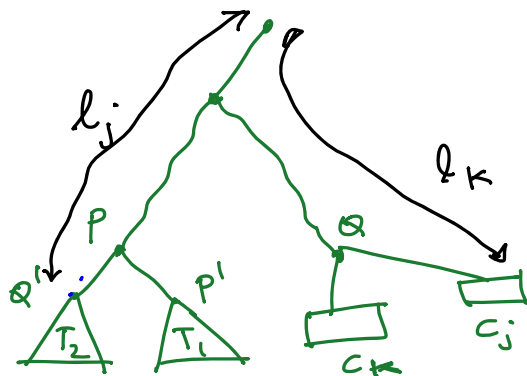
l_k :

.....
 pp' and qq' are single edges.

Let $l_k \geq l_j$ (the case $l_j \geq l_k$ is symmetrical)

switched the positions of c_j and T_2 .

T'



Let us assume the

$$wt(T_2) = \sum_{h=1}^r f_{i_h} \cdot d_{T_2}(c_{i_h})$$

$\{c_{i_1}, \dots, c_{i_r}\}$ are the leaf nodes in T_2

$$wt(T') - wt(T)$$

$$= \sum_{h=1}^r f_{i_h} (l_j + d_{T_2}(c_{i_h})) + f_j l_k - f_j l_j - \sum_{h=1}^r f_{i_h} (l_k + d_{T_2}(c_{i_h}))$$

$$= \underbrace{l_j \sum_{h=1}^r f_{i_h}}_d + \cancel{\sum_{h=1}^r f_{i_h} d_{T_2}(c_{i_h})} + f_j l_k - f_j l_j - \underbrace{l_k \sum_{h=1}^r f_{i_h}}_d - \cancel{\sum_{h=1}^r f_{i_h} d_{T_2}(c_{i_h})}$$

$$= l_j d + f_j l_k - f_j l_j - l_k d$$

$$= (l_j - l_k) d - (l_j - l_k) f_j$$

$$= (l_j - l_k) (d - f_j)$$

$$\leq 0 \quad \geq 0$$

$$\Rightarrow \leq 0$$

All the leaves in set $\{c_{i_1}, \dots, c_{i_r}\}$ have frequency

$$\geq f_j, f_k$$

$$\Rightarrow d - f_j \geq 0$$

$$\Rightarrow wt(T') \leq wt(T)$$

$$\Rightarrow wt(T') = wt(T)$$

($\because T$ is optimal)



Step 3

Claim:

There is an optimal solution to C which arises by combining greedy choice (for C) and optimal solution to the subproblem C' .

Proof:

By previous claim, there is an optimal solution S to C corresponding to greedy choice for C and solution S' to subproblem C' .

$$\underbrace{wt_C(S)}_{\text{optimal}} = \underbrace{wt_{C'}(S')}_{\text{Also optimal}} + (f_j + f_k)$$

where f_j, f_k are two least frequencies in C .

If S' was not optimal then there is S'' (a solution of C')
s.t. $wt_{C'}(S'') < wt_{C'}(S')$

\Rightarrow there is a solution R to C which arises by combining choice (f_j, f_k) with S'' .

$$wt_C(R) = wt_{C'}(S'') + (f_j + f_k)$$

$$< wt_{C'}(S') + (f_j + f_k) = wt_C(S)$$

A contradiction, because S is $\underset{\text{(min weight)}}{\text{optimal solution}}$ to C .

$\Rightarrow S'$ is optimal solution of C' .

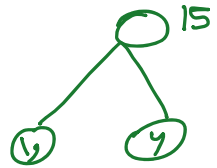


Example

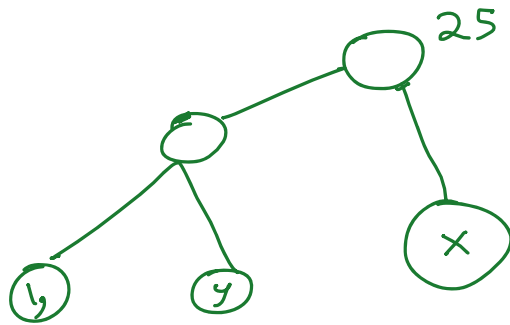
a, b, x, y, r, ~~l~~, l
30 25 10 5 20 20 10

← characters
← frequencies

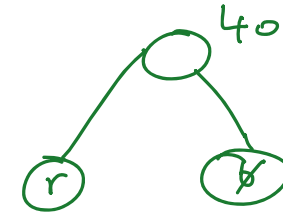
Compute optimal code



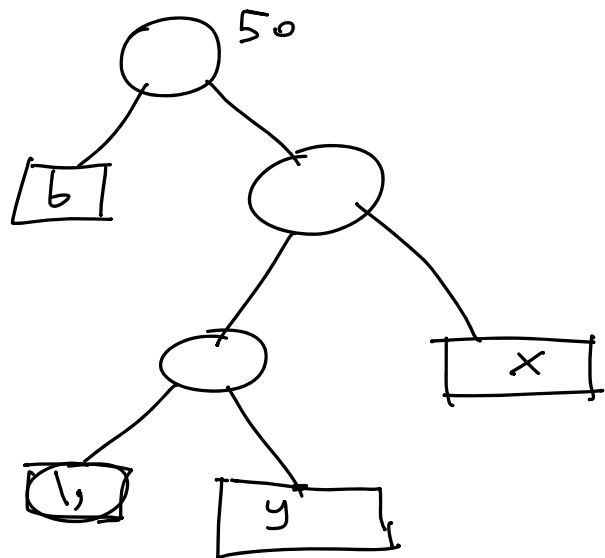
a b x r ~~l~~ l
30, 25 10 20 20 15



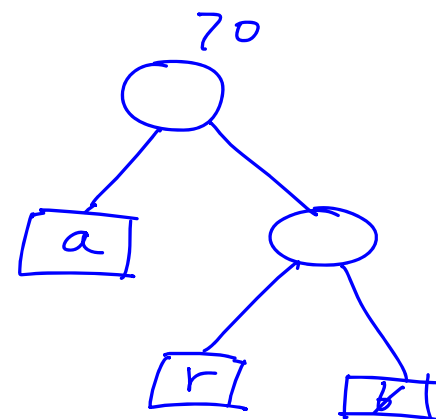
a b r ~~l~~ 25
30 25 20 20



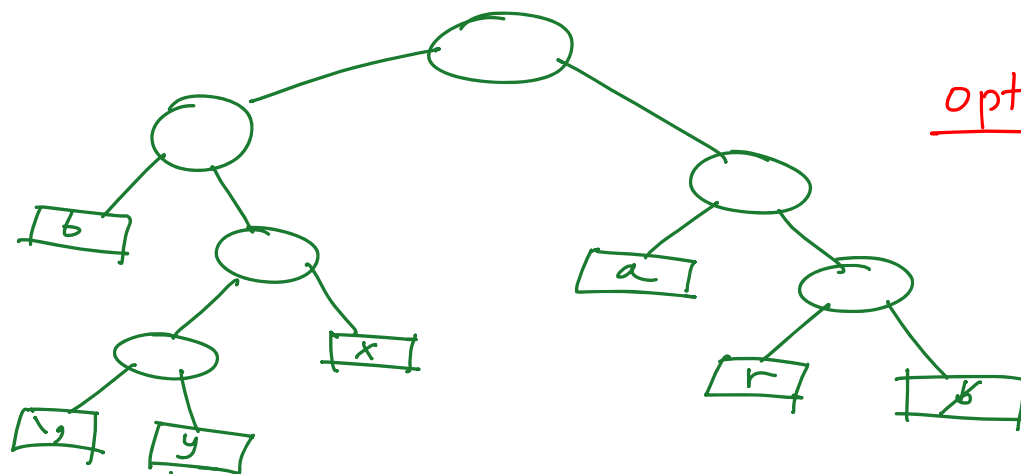
a b 40 25
30 25 25



$\begin{array}{r} a \\ \hline 30 \end{array} \quad \begin{array}{r} 40 \\ \hline \end{array} \quad 50$



50 70



optimal tree

$b \rightarrow 00$
 $a \rightarrow 10$
 $l \rightarrow 0100$
 $y \rightarrow 0101$
 $x \rightarrow 011$

Pseudo code

Huffman(C)

$n = |C|$

[Make priority queue Q of
element of C with freq being the key
for $i = 1$ to $n-1$ do

$x = \min(Q)$ // returns the min element

$\text{delete}(Q)$ // deletes the min element from Q

$y = \min(Q)$

$\text{delete}(Q)$

$z = \text{new node for tree}$

$z.lchild = x$

$z.rchild = y$

$z.freq = x.freq + y.freq$

$z.parent = \text{nil}$

$\text{insert}(Q, z)$

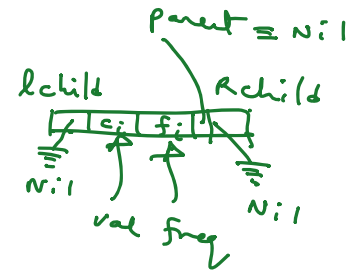
// end for

return $\min(Q)$

- At any stage of the algorithm we have a set of trees.
- Each tree has an associated freq.
- Each time we need to pick two trees with min weight

keep these trees in a priority queue.

Initially



Time Complexity $O(n \log n)$