

# COL 106 Lecture 14

Topic: Huffman Coding : Analysis

Announcement: Assignment 1 scores published

Recap:

Huffman Coding Problem

Input : Finite set  $S = \{a_1, \dots, a_n\}$  of "letters"  
A tve real number  $p_i$  for each  $a_i$ .

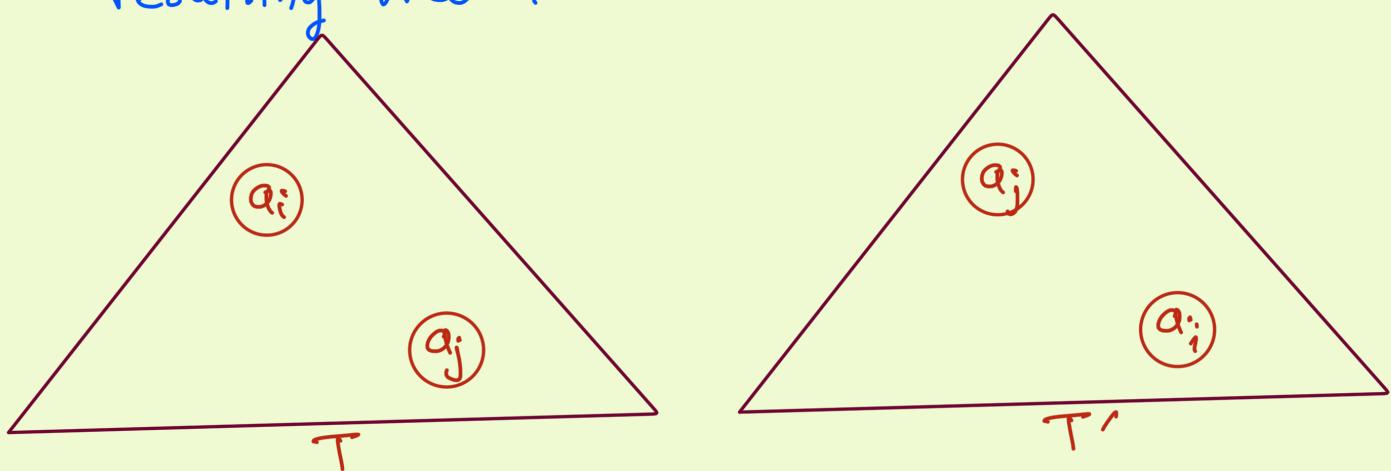
Output: A binary tree  $T$  with leaves  $a_1, \dots, a_n$   
that minimizes

$$\text{cost}(T) = \sum_{i=1}^n p_i \cdot \text{depth}(a_i).$$

Claim 1: If  $p_i < p_j$  then  $\text{depth}_T(a_i) \geq \text{depth}_T(a_j)$  in every optimum binary tree  $T$

Proof: By contradiction. Suppose  $p_i < p_j$  and  $\text{depth}_T(a_i) < \text{depth}_T(a_j)$  in the optimum binary tree, say  $T$ .

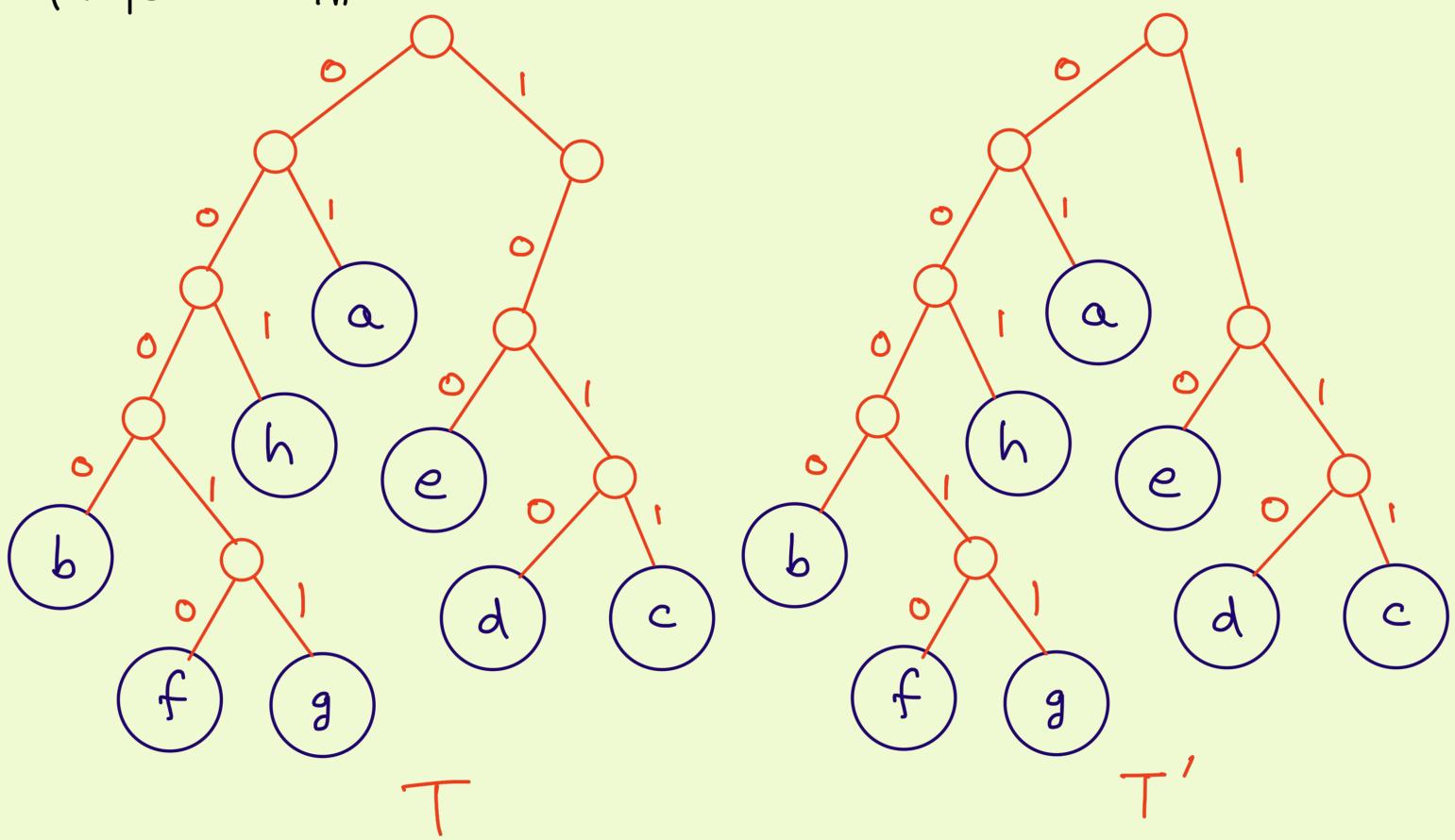
Exchange the locations of  $a_i$  and  $a_j$ . Call the resulting tree  $T'$ .



$$\begin{aligned}
 \text{cost}(T) - \text{cost}(T') &= p_i \underbrace{\text{depth}_T(a_i)}_{\text{depth } T(a_j)} + p_j \underbrace{\text{depth}_T(a_j)}_{\text{depth } T(a_i)} \\
 &\quad - p_i \underbrace{\text{depth}_{T'}(a_i)}_{\text{depth } T(a_j)} - p_j \underbrace{\text{depth}_{T'}(a_j)}_{\text{depth } T(a_i)} \\
 &= p_i (\underbrace{\text{depth}_T(a_i) - \text{depth}_{T'}(a_i)}_{< 0}) \\
 &\quad - p_j (\underbrace{\text{depth}_T(a_j) - \text{depth}_{T'}(a_j)}_{< 0}) \\
 &= (p_i - p_j) (\underbrace{\text{depth}_T(a_i) - \text{depth}_{T'}(a_j)}_{< 0}) > 0
 \end{aligned}$$

$\therefore T'$  is better than  $T \rightarrow$  contradiction to the assumption that  $T$  is an optimal tree.

$S = \{a, b, c, d, e, f, g, h\}$ . Can  $T$  be the optimal tree for some  $p_a, p_b, \dots, p_h > 0$ ?



Claim 2: An optimal tree can never contain a node having exactly 1 child.

Proof: Bypass such a node to get a better tree.

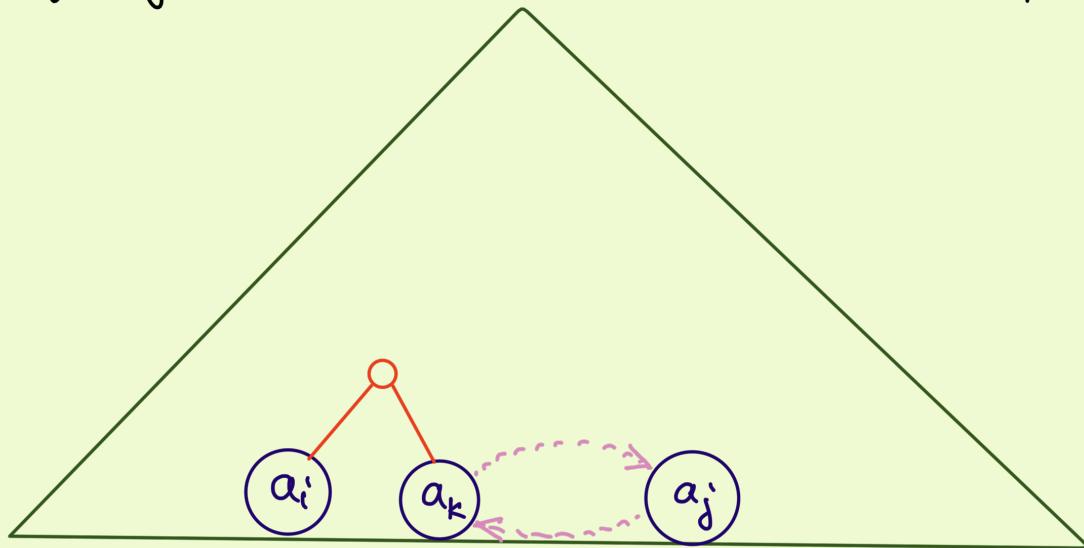
Let  $a_i$  : least frequent letter  $P_i \leq P_j \forall j$

Assume  $|S| \geq 2$

Let  $a_j$  : second least freq. letter.

Suppose  $a_k$  is the sibling of  $a_i$  and  $k \neq j$

Exchanging  $a_j$  and  $a_k$  gives us another optimal tree.



Claim: There exists an optimal tree in which the two least frequent letters are siblings. (Assuming  $n \geq 1$ )

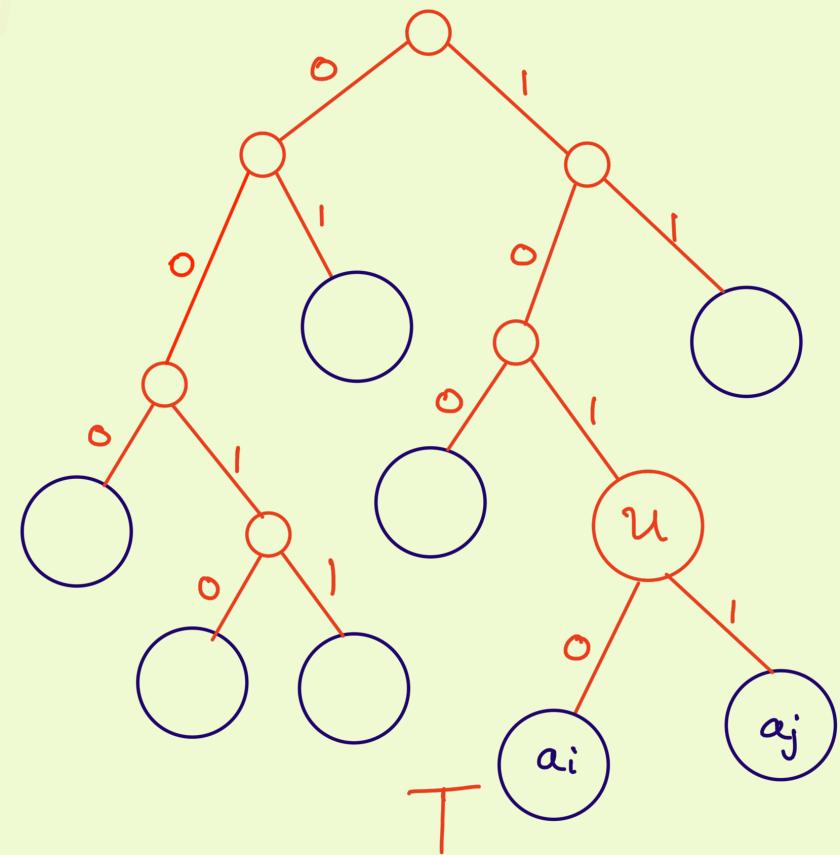
Proof:  $\geq 2$  leaves in the deepest level of an optimal tree.  
(due to claim 1)

$\therefore$  Two least frequent letters, say  $a_i, a_j$  in the bottom most level (due to claim 2)

$a_i$  has a sibling, say  $a_k$  (due to claim 1)

If  $k \neq j$ , exchange  $a_k$  and  $a_j$  to obtain another optimum tree.

$$\begin{aligned}
 \text{cost}(T) &= \sum_k p_k \cdot \text{depth}_T(a_k) \\
 &= \sum_{k \neq i, j} p_k \text{depth}_T(a_k) \\
 &\quad + p_i \text{depth}_T(a_i) \\
 &\quad + p_j \text{depth}_T(a_j) \\
 &= \sum_{k \neq i, j} p_k \text{depth}_T(a_k) \\
 &\quad + (p_i + p_j)(1 + \text{depth}(u)) \\
 &= \underbrace{\sum_{k \neq i, j} p_k \text{depth}_T(a_k)}_{\text{minimize}} + \underbrace{(p_i + p_j) \text{depth}(u)}_{\text{const.}} + p_i + p_j
 \end{aligned}$$



$$\begin{aligned}
 S = \{a_1, \dots, a_i, \dots, a_j, \dots, a_n\} &\rightarrow S' = (S \setminus \{a_i, a_j\}) \cup \{u\} \\
 p_1, \dots, p_i, \dots, p_j, \dots, p_n & \\
 p'_k &= p_k \quad \text{if } k \neq i, j \\
 p'_u &= p_i + p_j
 \end{aligned}$$

Recursive Algorithm:

1. Let  $a_i, a_j$  be two least freq. letters.
2. Construct a new instance with  
 $S' = (S \setminus \{a_i, a_j\}) \cup \{u\}$ ,  $p'_k = p_k$  for  $k \neq i, j$   
 $p'_u = p_i + p_j$
3. Recursively find the best tree, say  $T'$  for the new instance ( $T'$  has  $u$  as a leaf).
4. Attach  $a_i, a_j$  as children of  $u$ , and return this tree.