

- Huffman Tree

- Encoding Problem

- Given string S with n different character, length m

Find a way to encode S as a binary string C ,

Minimize the length of C .

- Fixed length encoding

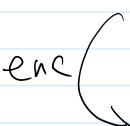
can use $T \lceil \log_2 n \rceil$ bits for every character.

Example: DNA Sequence $A=00$ $T=01$ $C=10$ $G=11$

$$\begin{array}{c} ATCCAG \longrightarrow \underline{\hspace{1cm}} \underline{\hspace{1cm}} \underline{\hspace{1cm}} \underline{\hspace{1cm}} \underline{\hspace{1cm}} \\ \text{m=6} \qquad \qquad \qquad A \ T \ C \ C \ AG \\ (\text{length} = m \cdot T \lceil \log_2 n \rceil) \end{array}$$

- Variable length encoding

idea: use a shorter length code for more frequent character.

ABAAAAAC

 $A = 0 \quad B = 10 \quad C = 11$
 $\begin{array}{cccccc} 0 & 10 & 00 & 00 & 11 & 0 \\ \hline A & B & A & A & A & C & A \end{array}$

problem: How to decode?

- Greedy Decoding character # $i \rightarrow a_i$ $\begin{pmatrix} a_1 = 0 \\ a_2 = 10 \\ a_3 = 11 \end{pmatrix}$

repeat

find shortest prefix of c such that $c[1..t] = a_i$

output i , remove $c[1..t]$

- want: a mapping $#i \rightarrow a_i$, so that greedy decoding can always correctly decode.

- "Prefix free": mapping $#i \rightarrow a_i$ is prefix-free
 if for any i, j a_i is not a prefix of a_j

$a_1 = 0$ prefix-free
 $a_2 = 10$
 $a_3 = 11$

$a_1 = 1$
 $a_2 = 10$
 $a_3 = 01$ not prefix-free

- Theorem: Greedy decoding works if and only if mapping is prefix-free.

Proof: if mapping is not prefix-free, $\exists a_i, a_j$

a_i is prefix of a_j

look at $\text{enc}(\#j)$, when you run greedy decoding
will output $\#\#i$, incorrect

if $c = \text{enc}(s)$, $s[i] = \#\#i$ $\text{len}(a_i) = t$

will prove: first output is $\#\#i$

$c[1..t] = a_i$
 $\forall t' < t \quad c[1..t'] = \text{Prefix of } a_i \neq a_j \quad (\text{prefix free})$
 \rightarrow greedy decoding will output $\#\#i$, correct.

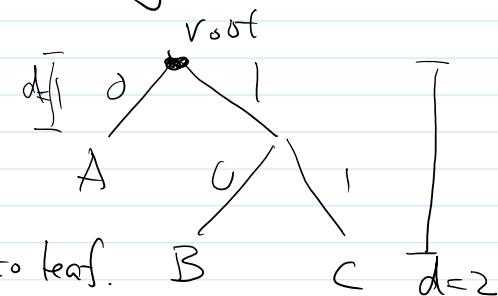
- Huffman tree: abstraction for prefix-free encoding

- Huffman tree is a rooted binary tree

- Each node has ≤ 2 children, labeled by 0, 1.

- Each leaf is labeled by a character

- Encoding of a character = path from root to leaf.



Claim: prefix-free mapping \iff Huffman tree

- Cost of Huffman Tree

$$w(T) = \sum_{i=1}^n d(\#i) \times f(\#i) = \text{total length of encoding.}$$

depth of character frequency of character $\#i$

Goal: minimize $w(T)$

Input: $n, f(\#i)$ for $i=1, 2, \dots, n$

repeat

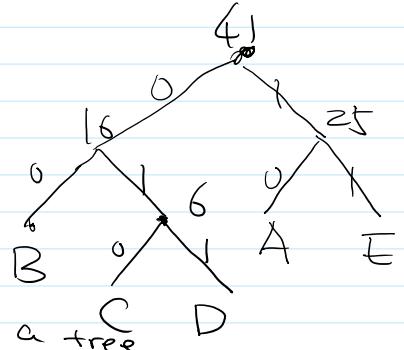
find two characters $\#i, \#j$ with smallest $f(\#i), f(\#j)$

merge them to a new character $\#new$

$$f(\#new) = f(\#i) + f(\#j)$$

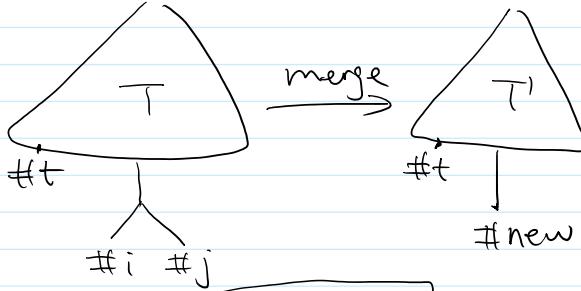
character A B CD E

frequency 12 10 5 1 13



- Theorem: Huffman tree algorithm always outputs a tree with smallest cost.

Observation



$$w(T) = w(T') + \boxed{f(\#i) + f(\#j)}$$

$$w(T) = \sum_{\{\#i, \#j\}} f(\#i) \times d(\#i) + f(\#j) \times d(\#j)$$

$$w(T') = \sum_{\{\#new\}} f(\#i) \times d(\#i) + f(\#new) \times d(\#new)$$

\uparrow
 $f(\#i) + f(\#j)$

$$\text{Cost of merging} = f(\#i) + f(\#j)$$

Claim: $w(T) = \sum_{\text{steps}} \text{cost of merging at this step}$

Proof (for optimality of Huffman Tree):

Do induction on n (# characters)

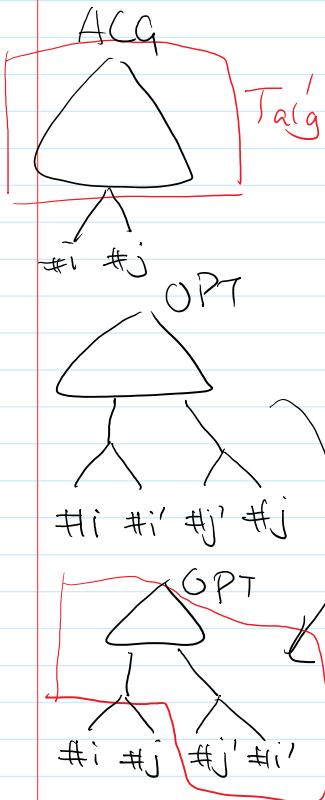
When $n=1$ there is only one tree.

induction

Suppose algorithm is optimal for $n \leq k$

Consider a case when $n=k+1$

Consider a case when $n = k+1$



Suppose in the first step ALG merged $\#i$, $\#j$

assume OPT is different: $\#i$, $\#j$ have different parents in OPT

if in OPT $d(\#i) \geq d(\#j)$

swap $\#j$ with $\#i'$

$$\begin{aligned} \text{new cost} &= \text{old cost} - d(\#i') \times f(\#i') - d(\#j) \times f(\#j) \\ &\quad + d(\#i') \times f(\#j) + d(\#j) \times f(\#i') \\ &= \text{old cost} - (f(\#i') - f(\#j)) \times (d(\#i') - d(\#j)) \\ &\geq 0 \quad \geq 0 \end{aligned}$$

(if $d(\#i) < d(\#j)$, swap $\#i$ with $\#j'$)

\Rightarrow There is an OPT with $\#i$, $\#j$ share the same parent.

ALG outputs optimal tree for all other characters + $\#i$, $\#j$

T_{alg} : algorithm's tree

T'_{alg} : alg's tree with $\#i$, $\#j$ merged

T_{OPT} : OPT 's tree

T'_{OPT} : OPT 's tree with $\#i$, $\#j$ merged.

$$w(T'_{alg}) = w(T'_{OPT}) \quad (\text{induction hypothesis})$$

$$w(T_{alg}) = w(T'_{alg}) + f(\#i) + f(\#j)$$

$$w(T_{OPT}) = w(T'_{OPT}) + f(\#i) + f(\#j)$$

$$w(T_{alg}) = w(T_{OPT})$$

□