

# Huffman Codes.

A-Z } 52  
a-z } 52  
0-9 } 10  
punctuation. } 10  
72  
62 + 10  
n

Quiz syllabus

→ Until MSTs

10:15  
11:15

With  $\lceil \log n \rceil$  many bits, we can represent  $n$  distinct letters.

→ English language.

e

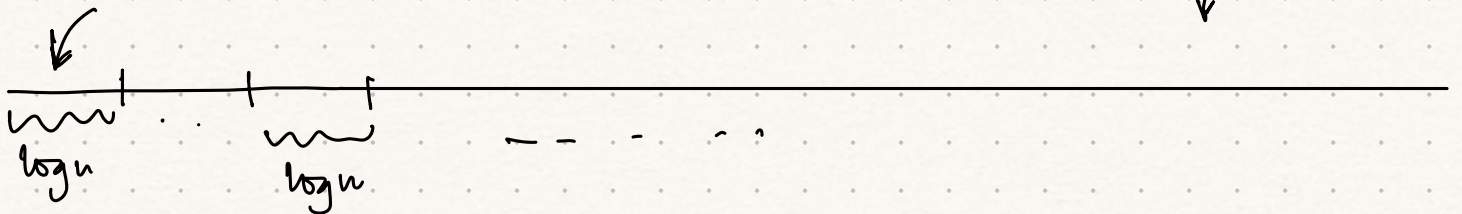
z, x

$(1\,000\,000) \times \lceil \log n \rceil$

Want an encoding s.t for any pair of letter  $u, v$ ,

$\left. \begin{array}{l} \text{Enc}(u) \text{ is not a prefix of } \text{Enc}(v) \\ \text{Enc}(v) \text{ " " } \text{Enc}(u) \end{array} \right\}$

Trivial Encoding using bin repr.



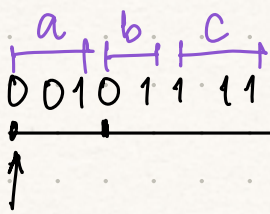
"Algorithm" →  $\text{Enc}(A) \cdot \text{Enc}(l) \cdot \text{Enc}(g) \dots \cdot \text{Enc}(m)$   
Binary string.

Here the length of  $\text{Enc}$  for all letters is the same.

Suppose, we move to a scheme where the lengths of  $\text{Enc}$  are different. Then, "not prefix" property helps.

$\left. \begin{array}{l} \text{Enc}(u) \text{ is not a prefix of } \text{Enc}(v) \\ \text{Enc}(v) \text{ " " } \text{Enc}(u) \end{array} \right\}$

a = 001  
b = 01  
c = 111



$d = 011$

Free ~ Prob

Letters	a	b	c	d	e
Freq	.32	.25	.2	.18	.05

In other words  
come up with a  
lookup table  
s.t ABL is  
optimized.

Task: Construct a scheme of Encoding s.t.

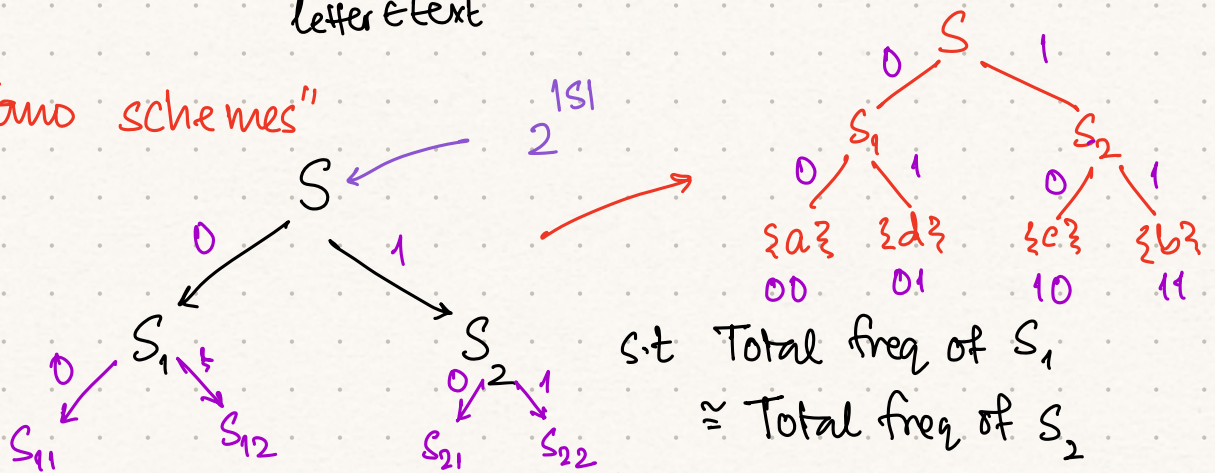
1. Optimize the average bit length (ABL)
2. "not prefix" property holds.

$$ABL(\text{text}) = \frac{\sum_{\text{letter} \in \text{text}} (\# \text{ of occurrences of letter}) \cdot |\text{Enc}(\text{letter})|}{\text{total length of text}}$$

Want to optimize  
ABL

$$= \sum_{\text{letter} \in \text{text}} \text{freq.}(\text{letter}) \cdot |\text{Enc}(\text{letter})|$$

"Shannon-Fano schemes"



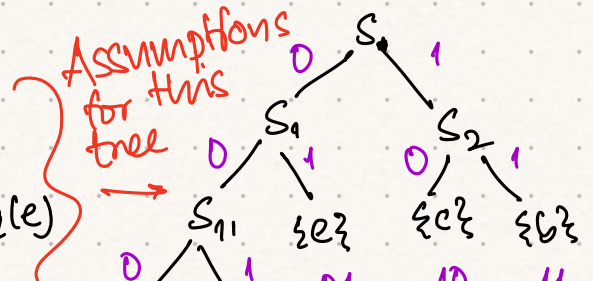
s.t Total freq of  $S_1$   
 $\approx$  Total freq of  $S_2$

Would give  $\log_2$  representation if we have a complete binary tree. That

- If we cannot find an equipartition.
- If there is a equipartition, searching for  $S_1$  and  $S_2$  could be tiresome.

$$\text{freq}(a) + \text{freq}(d) \approx \text{freq}(e)$$

$$\text{freq}(c) + \text{freq}(b) \approx \text{freq}(a) + \text{freq}(d) + \text{freq}(e)$$





$\text{freq}(a) \approx \text{freq}(d)$ .

$\{a\}$   $\{d\}$   
000 001

- We are labeling the leaves with letters and not the internal nodes.
  - Root to leaf path gives the Enc of the letter assoc. w/ the leaf.
- Ensure "not prefix" property.

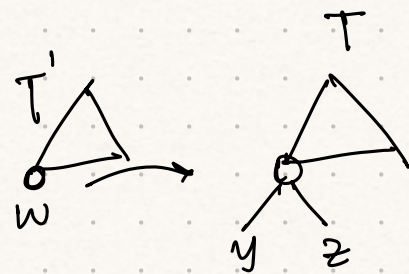
## Huffman encoding:

$S$ : Set of letters

$P$ : List of freq of letters in  $S$ .

Algo( $S, P$ ): Handle base cases

- Pick two least freq elements in  $S$  (say 'y' and 'z')
- Remove 'y' and 'z' from  $S$  and  $P$ .
- Add a new letter  $w$  to  $S$  s.t.  
 $\text{freq}(w) = \text{freq}(y) + \text{freq}(z)$



- Call updated lists  $S'$  and  $P'$
- Tree  $T' \leftarrow \text{Algo}(S', P')$

- Obtain a tree  $T$  by replacing the leaf  $w$  by a node with children 'y' and 'z'.

Return  $T$

$S$	a	b	c	d	e
$P$	.32	.25	.2	.18	.05

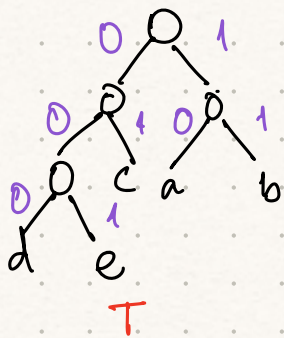
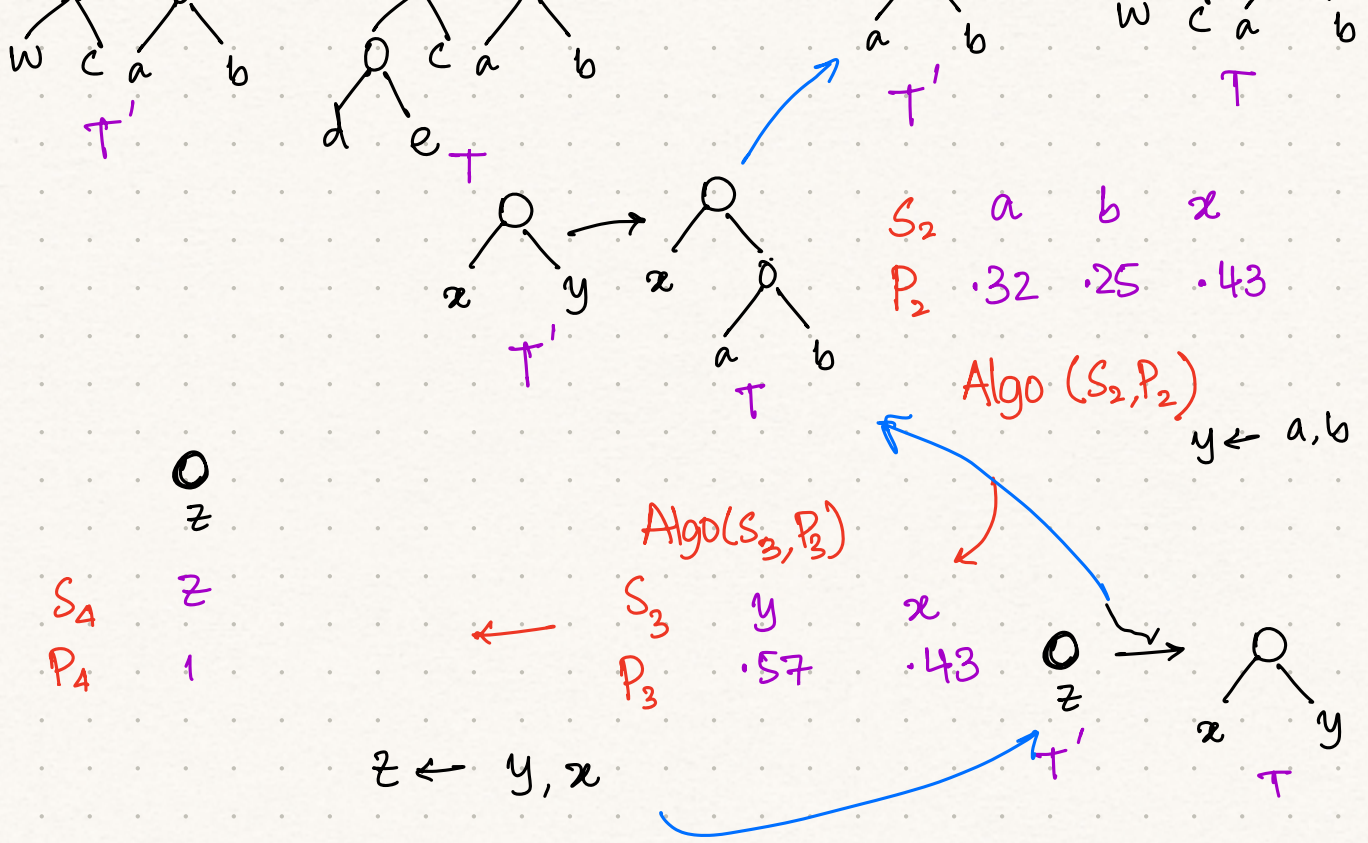
└──┬──┘

Algo( $S, P$ )

$S_1$	a	b	c	w
$P_1$	.32	.25	.2	.23

Algo( $S_1, P_1$ )





a — 10  
 b — 11  
 c — 01  
 d — 000  
 e — 001

$S$  a b c d e  
 $P$  .32 .25 .2 .18 .05

$$\begin{aligned}
 & 2 \times .32 \\
 & + 2 \times .25 \\
 & + 2 \times .2 \\
 & + 3 \times .18 \\
 & + 3 \times .05
 \end{aligned}
 \left. \vphantom{\begin{aligned} & 2 \times .32 \\ & + 2 \times .25 \\ & + 2 \times .2 \\ & + 3 \times .18 \\ & + 3 \times .05 \end{aligned}} \right\} = 2.23$$

ABL<sub>T</sub>  
 Better than 3

Huffman encoding is optimal

Alphabet  
 $S \leftarrow$  Set of letters  
 $P \leftarrow$  freq list.

→ The tree and the corresponding encoding is optimal.

$$ABL_{\tilde{T}^*} \leq ABL_T \quad \forall T \text{ that follows "not prefix" property}$$

• Proof by induction on  $|S|$ .



Base case:  $|S|=1$  Trivial.

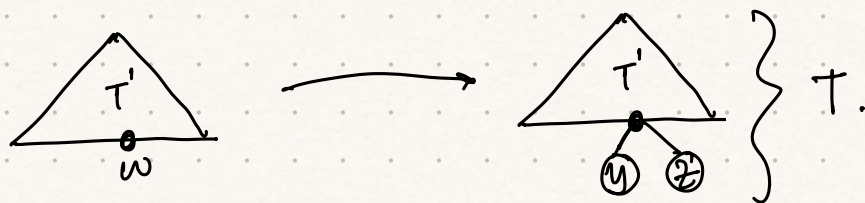
I.H: True for  $|S| \leq k-1$ .

Induction Step: For the sake of contradiction, Huffman encoding gives a tree  $\underline{T}$  but  $\exists Z$  s.t.  $ABL_Z < ABL_T$ . — (\*)

- Picks 2 least freq elements  $y$  and  $z$  and replaces them in  $S$  with a new elem ' $w$ '.

$$S' \leftarrow (S \setminus \{y, z\}) \cup \{w\}. \quad |S'| = k-1.$$

- From ind. hyp., we get an optimal tree  $T'$ .



$$\text{depth}_T(y) = \text{depth}_T(z)$$

$$\text{freq}(w) = \text{freq}(y) + \text{freq}(z)$$

$$ABL_T(S) = \sum_{u \in S} \text{freq}(u) \cdot \text{depth}_T(u)$$

$$= \sum_{\substack{u \in S \setminus \{y, z\} \\ \underbrace{\hspace{1cm}} \\ S' \setminus \{w\}}} \text{freq}(u) \cdot \text{depth}_T(u) + (\text{freq}(y) + \text{freq}(z)) \cdot \text{depth}_T(y).$$

$$= \sum_{u \in S' \setminus \{w\}} \text{freq}(u) \cdot \text{depth}_T(u) + \text{freq}(w) \cdot (\text{depth}_T(w) + 1)$$

$$= \left( \sum_{u \in S'} \text{freq}(u) \cdot \text{depth}_T(u) \right) + \text{freq}(w)$$

$$= ABL_{T'}(S') + \text{freq}(w) \quad \text{--- (1)}$$

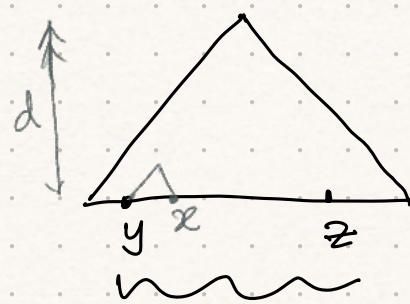
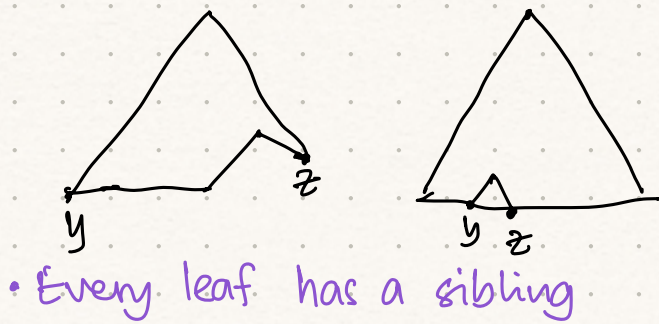
$$\bullet \forall u \in S \setminus \{y, z\}, \text{depth}_T(u)$$

$$= \text{depth}_{T'}(u)$$

$$\bullet \text{depth}_T(y) = \text{depth}_{T'}(w) + 1$$

Let us consider  $Z$ .

Qn: Where do you think 'y' and 'z' are, in tree  $Z$ ?



Subtle exchange / swap argument.

$$\text{Rest} + \text{freq}(x) \cdot d + \text{freq}(y) \cdot d + \text{freq}(z) \cdot d$$

W.L.O.G., assume that 'y' and 'z' are siblings in the tree  $Z$ .

Obtain  $Z'$  by replacing parent of 'y' and 'z' with  $w$ .

$$ABL_Z(S) = ABL_{Z'}(S') + \text{freq}(w) \quad \text{--- (2)}$$

$$ABL_{Z'}(S') + \text{freq}(w) = ABL_Z(S) < ABL_T(S) = ABL_{T'}(S') + \text{freq}(w)$$

$\Downarrow$

$$ABL_{Z'}(S') < ABL_{T'}(S').$$

Contradicts the optimality of  $T'$  w.r.t set  $S'$  obtained through I.H.