

## Huffman code - Proof of optimality:

Claim of optimality:

Let  $\bar{L}_{\text{Huffman}}$  be the length of a code constructed for  $X \sim (p_1, \dots, p_k)$  according to the Huffman algorithm. Then  $\bar{L}_{\text{Huffman}} = \min_{C \in \text{Prefix free codes for } X} \bar{L}_C$ .

Proof:

We will denote the avg length of some code  $\ell$  for some random variable taking  $M$  values with probabilities  $q_1, \dots, q_M$  as  $L_{\ell}(q_1, \dots, q_M)$ .

The optimal length for the same distribution is written as  $L^*(q_1, \dots, q_m) \rightarrow$  [note that this <sup>number</sup> does not depend on which optimal code we are choosing]

We will show that

$L^*(p_1, \dots, p_k) = L^*(p_1, \dots, p_{k-1} + p_k) + p_{k-1} + p_k$   
 2 least possibility symbols are combined

& then this will imply optimality of Huffman.

2 least probability symbols  
are combined

$\hookrightarrow A$

↙  
To  
be  
~~proved~~  
later

How?  $\swarrow$   
 $\underline{\underline{=}}$

An optimal code for  $(p_1, \dots, p_k)$  can be obtained from an optimal code for  $(p_1, \dots, p_{k-1} + p_k)$  in the following way

$l_i$

$(l'_{k-1})$   
 $\leftarrow$   $p_{k-1} + p_k$   
 " length of this leaf from root

Rest of binary tree of code with length  $L^*$   
 $(p_1, \dots, p_{k-1} + p_k)$

Consider a code  $\mathcal{C}$  for  $(p_1, \dots, p_k)$

$(l'_{k-1} + 1)$   
 $\leftarrow$  (leaf) Codeword for  $p_{k-1}$   
 $(l'_{k-1} + 1)$   
 $\leftarrow$  (leaf) Codeword for  $p_k$

No change

Avg length of  $\mathcal{C} = \bar{L}_{\mathcal{C}} = \sum_{i=1}^k p_i l_i$

$$= \sum_{i=1}^{k-2} p_i l'_i + p_{k-1} (l'_{k-1} + 1) + p_k (l'_{k-1} + 1)$$

$$= \left( \bar{L}^*_{(p_1, \dots, p_{k-1} + p_k)} + (p_{k-1} + p_k) \right)$$

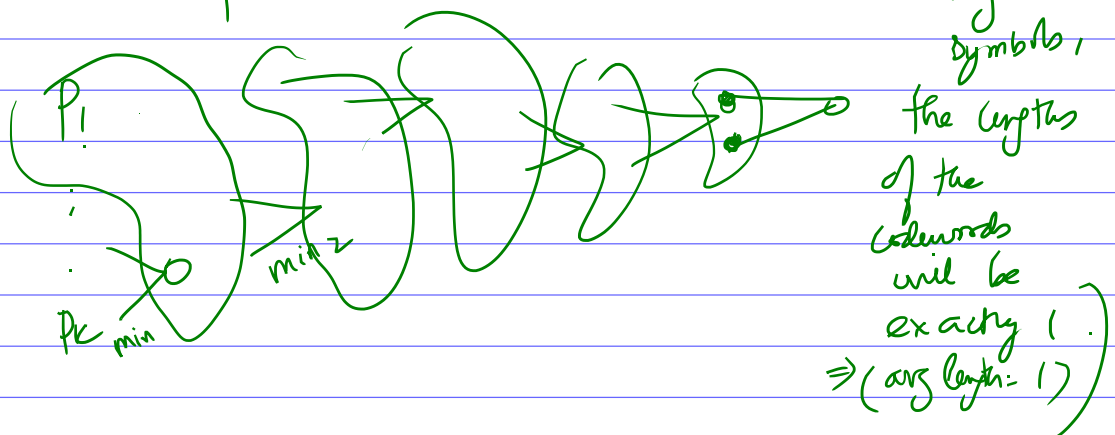
This is optimal by (A) for  $(p_1, \dots, p_k)$

$$(C) = L^*(p_1, \dots, p_k)$$

To get a optimal code for  $(p_1, \dots, p_{k-1} + p_k)$

We obtain an optimal code for  $(k-2)$  - prob distribution obtained by combining 2 least prob symbols of  $(p_1, \dots, p_{k-1} + p_k)$

& keep repeating this process until we have only 2 probabilities in the distribution. (When there are only 2 symbols,



Huffman algorithm constructs such a binary tree

& assigns codewords such that at each stage the

$$\rightarrow \text{avg length at this stage} = \left( \text{sum of smallest 2 probabilities} \right) + \text{avg length of code at next stage}$$

$\rightarrow$  So optimal length is assured by above formula since in the last stage, Huffman algorithm exactly gets an (with 2 prob) optimal code

$$\text{This shows that } (\text{Prop A}) \Rightarrow L_{\text{Huffman}} = L_{(p_1, p_k)}^*$$

Now we show Prop A: If  $k=2$  (then the statement of Prop A is trivial as  $L_{k-1}^* = 0$ )

$k > 2$  Let  $c_k^*$  be some optimal code for  $(p_1, \dots, p_k)$

&  $c_{k-1}^*$  be optimal code for  $(p_1, \dots, p_{k-1} + p_k)$

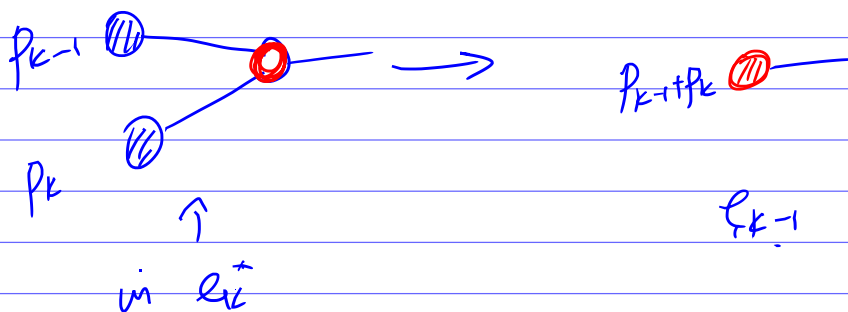
$$L_{c_k^*} = L_{(p_1, \dots, p_k)}^*$$

$$L_{c_{k-1}^*} = L_{(p_1, \dots, p_{k-1} + p_k)}^*$$

Step 1: Now from  $\mathcal{C}_k^*$  we will obtain a new code for  $\mathcal{C}_{k-1}$   
 $(p_1, \dots, p_{k-1} + p_k)$  with length  $\bar{L}_{\mathcal{C}_{k-1}} = \bar{L}_{(p_1, \dots, p_k)}^* - p_{k-1} - p_k$

Step 2: also, from  $\mathcal{C}_{k-1}^*$ , we will obtain a new code  $\mathcal{C}_k$   
for  $(p_1, \dots, p_k)$  with length  $\bar{L}_{\mathcal{C}_k} = \bar{L}_{(p_1, \dots, p_{k-1} + p_k)}^* + p_{k-1} + p_k$

Done Step 1: We assume that the code  $\mathcal{C}_k^*$  is picked so that lemma 3 is satisfied



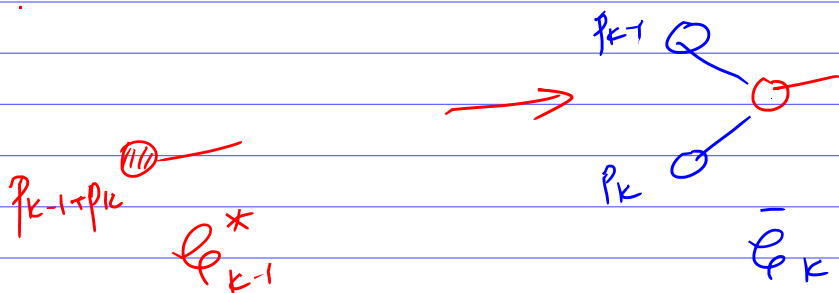
From the above construction, we see that

$$\bar{L}_{\mathcal{C}_{k-1}} = \bar{L}_{\mathcal{C}_k^*} - p_{k-1} - p_k$$

$$\Rightarrow \bar{L}_{\mathcal{C}_k^*} - \bar{L}_{\mathcal{C}_{k-1}} = p_{k-1} + p_k$$

$$\Rightarrow \bar{L}_{(p_1, \dots, p_k)}^* - \bar{L}_{\mathcal{C}_{k-1}} = p_{k-1} + p_k \rightarrow \textcircled{1}$$

Done Step 2:



$$\bar{L}_{\mathcal{C}_k} = \bar{L}_{\mathcal{C}_{k-1}^*} + p_{k-1} + p_k$$

$$\bar{L}_{\mathcal{C}_k} = \bar{L}_{(p_1, \dots, p_{k-1} + p_k)}^* + p_{k-1} + p_k$$

$$\bar{L}_{q_k} - \bar{L}^*_{(p_1, \dots, p_{k-1} + p_k)} = p_{k-1} + p_k \rightarrow (2)$$

By ① & ② as RHS is same, LHS must be same

$$\Rightarrow \bar{L}_{q_k} - \bar{L}^*_{(p_1, \dots, p_{k-1} + p_k)} = \bar{L}^*_{p_1, \dots, p_k} - \bar{L}_{q_{k-1}}$$

$$\Rightarrow \bar{L}_{q_k} - \bar{L}^*_{(p_1, \dots, p_k)} = \bar{L}^*_{(p_1, \dots, p_{k-1} + p_k)} - \bar{L}_{q_{k-1}}$$

$\downarrow$   
 $\geq 0$

$\downarrow$   
 $\leq 0$

$\Rightarrow$  both sides = 0

$$\Rightarrow \bar{L}_{q_k} = \bar{L}^*_{p_1, \dots, p_k} \rightarrow (3)$$

$$\& \bar{L}^*_{p_1, \dots, p_{k-1} + p_k} = \bar{L}_{q_{k-1}} \rightarrow (4)$$

$\Rightarrow$  Prop A is true (Take (3) & substitute in (2))

This completes the proof of optimality of Huffman coding algorithm.