

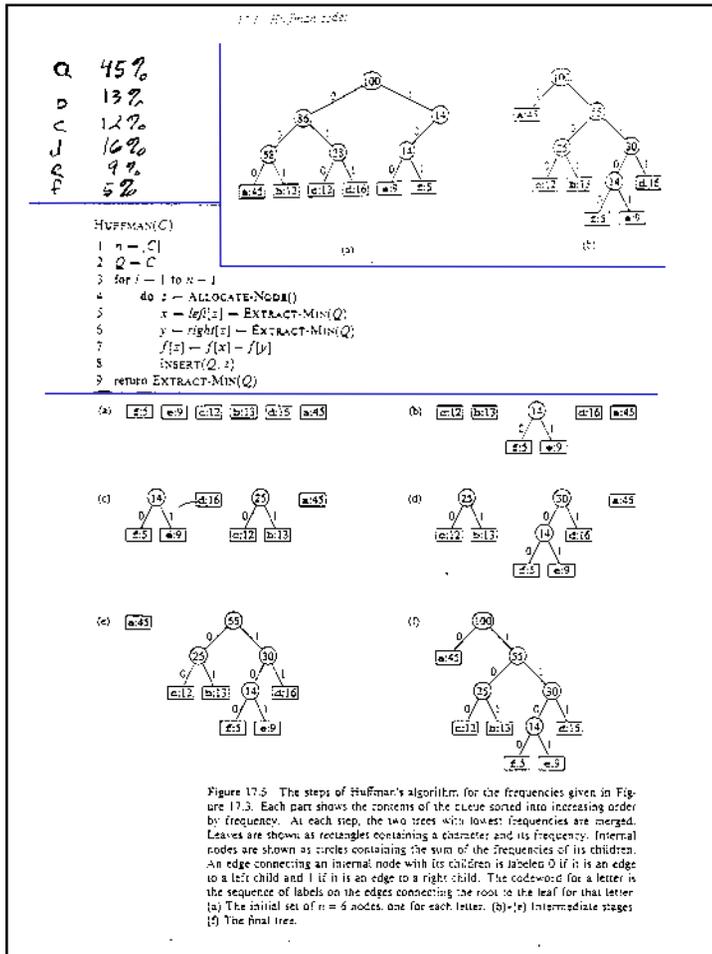
# Huffman Codes

---

An Optimal Data  
Compression Method

## Data Compression

- Binary character code (“code”)
  - each k-bit source string maps to unique code word (e.g. k=8)
  - “compression” alg: concatenate code words for successive k-bit “characters” of source
- Fixed/variable length codes
  - all code words equal length?
- Prefix codes
  - no code word is prefix of another (simplifies decoding)

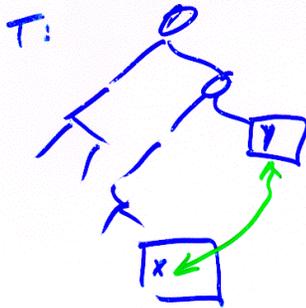


## Strategy

1. Opt solution maybe not unique  
(So cannot prove greedy gives only possible answer.)
2. Show greedy's solution is as good as any.

## Huffman Code Objective

$$\begin{aligned}
 \text{compressed length} &= \\
 &= \sum_{c \in C} \text{freq}(c) \cdot \text{length}(c) \\
 &= \sum_{c \in C} \text{freq}(c) \cdot \text{depth}(c) \\
 &= B(T)
 \end{aligned}$$



Defn: a pair of leaves  $x, y$  is an inversion if  
 $\text{depth}(x) \geq \text{depth}(y)$   
 $\& \text{freq}(x) \geq \text{freq}(y)$

claim if we flip an inversion cost never increases

Why?: all other things being equal, better to give more frequent letter shorter code

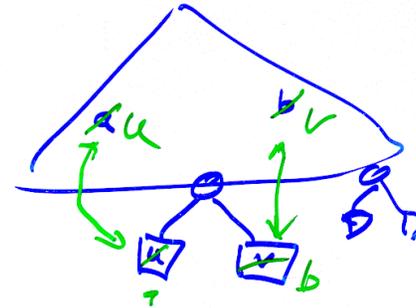
$$\underbrace{d(x)f(x) + d(y)f(y)}_{\text{before}} - \underbrace{(d(x)f(y) + d(y)f(x))}_{\text{after}}$$

$$= (d(x) - d(y)) (f(x) - f(y)) \geq 0$$

i.e. cost savings is non-negative

Lemma 17.2

2 least frequent chars might as well be siblings, at deepest level.



"greedy choice property"

- Let  $a$  be least freq letter,  $b$  2<sup>nd</sup>.
- Let  $u$  be least freq leaf at max depth,  $v$  its sibling.
- Then  $(a, u)$  &  $(b, v)$  are inversions. Swap them.

Modified lemma 17.3

"optimal substructure"

Let  $\mathcal{C}$  be an  $n$ -letter alphabet with frequencies  $f(c)$ .

$\forall x, y \in \mathcal{C}$ , let  $\mathcal{C}'$  be  $(n-1)$ -letter alphabet  $\mathcal{C} - \{x, y\} \cup \{z\}$  with

$$\text{freq } f'(c) = \begin{cases} f(c), & c \neq x, y, z \\ f(x) + f(y), & c = z \end{cases}$$

Let  $T'$  be optimal tree for  $\mathcal{C}'$ .

Then  $T = T' - z + \begin{matrix} \circ \\ \swarrow \searrow \\ \square \square \end{matrix}$  is opt for  $\mathcal{C}$  among trees having  $x, y$  as siblings

Proof:

$$B(T) = \sum_{c \in \mathcal{C}} d_T(c) \cdot f(c)$$

$$B(T) - B(T') = d_T(x)(f(x) + f(y)) - d_{T'}(z) f'(z)$$

$$= (d_{T'}(z) + 1) f'(z)$$

$$- d_{T'}(z) f'(z)$$

$$= f'(z)$$

Suppose  $\hat{T}$  <sup>(having  $x, y$  as siblings)</sup> better than  $T$ :  $B(\hat{T}) < B(T)$   
Collapse  $x, y$  to  $z$ , forming  $\hat{T}'$ .

$$B(\hat{T}) - B(\hat{T}') = f'(z)$$

Then

$$B(\hat{T}') = B(\hat{T}) - f'(z) < B(T) - f'(z) = B(T')$$

contradiction

### Theorem 17.4

Huffman alg. gives optimal code.

Proof: by induction on # of letters.

basis:  $n=1$  trivial

ind:  $n \geq 2$

- let  $x, y$  be least frequent.
- Form  $\mathbb{Z}, \mathbb{C}'$  as above
- By induction  $T'$  is opt for  $\mathbb{C}'$
- By <sup>modified</sup> 17.3,  $T' \rightarrow T$  is opt for  $\mathbb{C}$  among trees with  $x, y$  siblings
- By 17.2 some opt. tree does have  $x, y$  siblings
- ∴  $T$  is optimal.

### Data Compression

Huffman is optimal.

But still might do better!

- ① Huffman encodes fixed length blocks.
- ② Huffman uses one encoding throughout a file. What if characteristics change?
- ③ What if data has structure?  
Eg. raster images  
Video
- ④ Lossless  
Ziv-Lempel, MPEG, ...