

# CSE 421 Intro to Algorithms Winter 2004

## Huffman Codes: An Optimal Data Compression Method

1

## Data Compression

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

CSE 421, W'04, Ruzzo

2

## Compression Example

- n 100k file, 6 letter alphabet:
- n File Size:
  - n ASCII, 8 bits/char: 800kbits
  - n  $2^3 > 6$ ; 3 bits/char: 300kbits
  - n 00,01,10 for a,b,d; 11xx for c,e,f: 2.52 bits/char  $74\% \cdot 2 + 26\% \cdot 4$ : 252kbits
  - n Optimal?
- n Why?
  - n Storage, transmission vs 1Ghz cpu

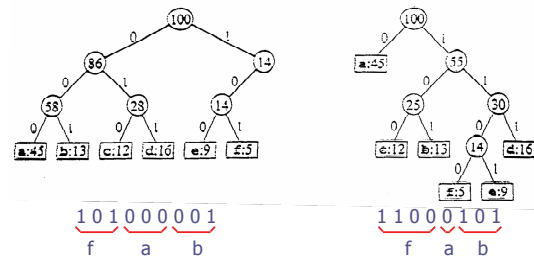
a	45%
b	13%
c	12%
d	16%
e	9%
f	5%

CSE 421, W'04, Ruzzo

3

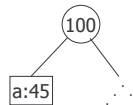
## Prefix Codes = Trees

a	45%
b	13%
c	12%
d	16%
e	9%
f	5%



## Greedy Idea #1

- n Put most frequent under root, then recurse ...



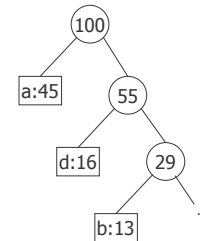
CSE 421, W'04, Ruzzo

5

## Greedy Idea #1

- n Put most frequent under root, then recurse ...

**Too greedy:  
unbalanced tree**



CSE 421, W'04, Ruzzo

6

## Greedy idea #2

a	45%
b	13%
c	12%
d	16%
e	9%
f	5%

Group least frequent letters near bottom

CSE 421, W'04, Ruzzo 7

CSE 421, W'04, Ruzzo 10

CSE 421, W'04, Ruzzo 10

## Huffman's Algorithm (1952)

**Algorithm:**

- insert node for each letter into priority queue by freq
- while queue length > 1 do
- remove smallest 2; call them x, y
- make new node z from them, with  $f(z) = f(x) + f(y)$
- insert z into queue

**Analysis:**  $O(n)$  heap ops:  $O(n \log n)$

**Goal:** Minimize  $B(T) = \sum_{c \in C} \text{freq}(c) * \text{depth}(c)$

**Correctness:** ???

CSE 421, W'04, Ruzzo 10

## Correctness Strategy

Optimal solution may not be unique, so cannot prove that greedy gives the only possible answer.

Instead, show that greedy's solution is as good as any.

CSE 421, W'04, Ruzzo 11

**Defn:** A pair of leaves is an inversion if

- $\text{depth}(x) \geq \text{depth}(y)$
- and
- $\text{freq}(x) \geq \text{freq}(y)$

**Claim:** If we flip an inversion, cost never increases.

**Why?** All other things being equal, better to give more frequent letter the shorter code.

**before**                      **after**

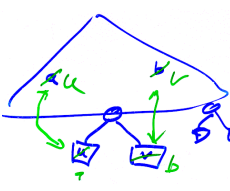
$$(d(x) * f(x) + d(y) * f(y)) - (d(x) * f(y) + d(y) * f(x)) = (d(x) - d(y)) * (f(x) - f(y)) \geq 0$$

i.e. non-negative cost savings.

### Lemma 1: “Greedy Choice Property”

The 2 least frequent letters might as well be siblings at deepest level

- Let  $a$  be least freq,  $b$   $2^{\text{nd}}$
- Let  $u, v$  be siblings at max depth,  $f(u) \leq f(v)$
- Then  $(a, u)$  and  $(b, v)$  are inversions. Swap them.



CSE 421, W'04, Ruzzo 13

### Lemma 2: “Optimal Substructure”

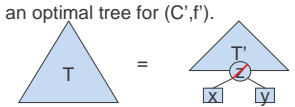
Let  $(C, f)$  be a problem instance:  $C$  an  $n$ -letter alphabet with letter frequencies  $f(c)$  for  $c$  in  $C$ .

For any  $x, y$  in  $C$ , let  $C'$  be the  $(n-1)$  letter alphabet  $C - \{x, y\} \cup \{z\}$  and for all  $c$  in  $C'$  define

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

Let  $T'$  be an optimal tree for  $(C', f')$ .

Then



is optimal for  $(C, f)$  among all trees having  $x, y$  as siblings

CSE 421, W'04, Ruzzo 14

Proof:

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

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

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

$$= f'(z)$$

Suppose  $\hat{T}$  (having  $x$  &  $y$  as siblings) is better than  $T$ , i.e.  $B(\hat{T}) < B(T)$ . Collapse  $x$  &  $y$  to  $z$ , forming  $\hat{T}'$

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

Then:

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

Contradicting optimality of  $T'$

### Theorem: Huffman gives optimal codes

Proof: induction on  $|C|$

- Basis:  $n=1, 2$  – immediate
- Induction:  $n > 2$ 
  - Let  $x, y$  be least frequent
  - Form  $C', f'$ , &  $z$ , as above
  - By induction,  $T'$  is opt for  $(C', f')$
  - By lemma 2,  $T' \rightarrow T$  is opt for  $(C, f)$  among trees with  $x, y$  as siblings
  - By lemma 1, some opt tree has  $x, y$  as siblings
  - Therefore,  $T$  is optimal.

CSE 421, W'04, Ruzzo 16

### Data Compression

- Huffman is **optimal**.
- BUT** still might do better!
  - Huffman encodes fixed length blocks. What if we vary them?
  - Huffman uses one encoding throughout a file. What if characteristics change?
  - What if data has structure? E.g. raster images, video, ...
  - Huffman is lossless. Necessary?
- LZW, MPEG, ...

CSE 421, W'04, Ruzzo 17