CMPT307: Greedy Algorithms

Week 8-2

Xian Qiu

Simon Fraser University



Encoding Characters

to encode 100 characters: a-f only

using binary

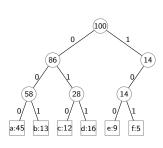
	a	b	С	d	е	f
frequency	45	13	12	16	9	5
fixed length codeword	000	001	010	011	100	101
variable length codeword	0	101	100	111	1101	1100

⊳ fixed length: 300 bits

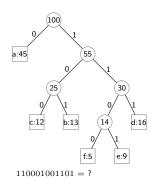
variable length: 224 bits

□ consider prefix code: no codeword is a prefix of another

Decoding



characters are leaves



$$f_c = \text{ frequency}, \ d_T(c) = \text{ depth}, \quad \forall c \in C = \text{alphabet}$$
 total bits:
$$B(T) = \sum_{c \in C} f_c \cdot d_T(c)$$

Coding Problem

- \triangleright input: alphabet C, with frequency $f_c > 0$ for all $c \in C$
- ightharpoonup output: prefix codes, minimizing the total bits equivalently, represent by a binary tree T minimizing B(T)

Definition

A binary tree is full if each internal node has exactly two children.



 \triangleright denote by T^* the optimal binary tree

Properties of OPT

Lemma 1. T^* is full.

Proof

assume T^* is not full



 T^{\prime} is a feasible encoding and has less bits than T^{*}

Properties of OPT

Lemma 2

 T^* has |C| leaves and |C|-1 internal nodes.

- ightharpoonup removing all leaves of T^* yields a new tree T_1^* with $\frac{|C|}{2}$ leaves
- ightarrow removing all leaves of T_1^* yields T_2^* with $\frac{|C|}{2^2}$ leaves
- ho continue this way until T_h^* with $\frac{|C|}{2^h}=1$ leave (singleton)
- b total number of internal nodes is

$$\begin{aligned} \frac{|C|}{2} + \frac{|C|}{2^2} + \dots + \frac{|C|}{2^h} &= |C| \frac{\frac{1}{2} \left(1 - \frac{1}{2^h}\right)}{1 - \frac{1}{2}} \\ &= |C| \left(1 - \frac{1}{|C|}\right) = |C| - 1 \end{aligned}$$

Problem Reformulation

- \triangleright given n nodes c_1,\ldots,c_n and their frequencies f_{c_1},\ldots,f_{c_n}
- \triangleright construct a full binary tree T^* of 2n-1 nodes s.t.
 - 1. c_1, \ldots, c_n are all leaves of T^*
 - 2. $B(T^*) = \sum_{i=1}^n f_{c_i} \cdot d_{T^*}(c_i)$ is minimized

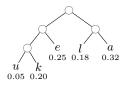
greedy strategies

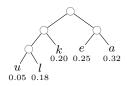
- \triangleright since T^* is full, consider two nodes at a time
- ▷ top-down: letters with high frequency should be near the top
- ▷ bottom-up: lowest frequency letters should be at the lowest level

Top-Down

algorithm: divide C into sets C_1 and C_2 with (almost) equal frequencies, and recursively build tree for C_1 and C_2

Example. $f_a = 0.32, f_e = 0.25, f_k = 0.20, f_l = 0.18, f_u = 0.05.$





the algorithm does not yield optimal solution!

Bottom-Up

observations

- \triangleright for n > 1, the lowest level always contains at least two leaves
- by the order in which items appear in a level does not matter

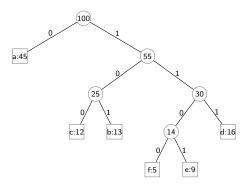
Lemma

Assume $x, y \in C$ have lowest frequencies, then x and y are siblings of T^* .

 \triangleright note that x and y are leaves

algorithm: make two leaves for x and y, recursively build tree for the rest using a meta-letter "xy"

Example



how to find the min-frequency node?

Pseudocode

```
c.f = frequency of c \in C
```

Q is min-priority queue keyed by c.f for all $c \in C$

Huffman(C)

1 n = |C|; 2 Q = C;

```
3 for i=1 to n-1 do

4 allocate a new node z;

5 z.left = x = \text{EXTRACT-MIN}(Q);

6 z.right = y = \text{EXTRACT-MIN}(Q);

7 z.freq = x.f + y.f;

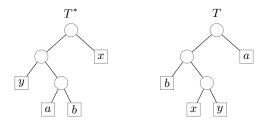
8 INSERT(Q, z);

9 return EXTRACT-MIN(Q); // return the root of the tree
```

Correctness

Lemma

Assume $x, y \in C$ have lowest frequencies, then x and y are siblings of T^* .



- $\triangleright a, b$ are siblings of maximum depth
- \triangleright assume $a \neq x$ and $y \neq b$ (what if not?)

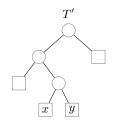
Correctness

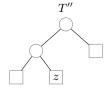
$$\begin{split} B(T) - B(T^*) &= \sum_{c \in C} c. \mathbf{f} \cdot d_T(c) - \sum_{c \in C} c. \mathbf{f} \cdot d_{T^*}(c) \\ &= x. \mathbf{f} \cdot d_T(x) + a. \mathbf{f} \cdot d_T(a) - x. \mathbf{f} \cdot d_{T^*}(x) - a. \mathbf{f} \cdot d_{T^*}(a) \\ &+ y. \mathbf{f} \cdot d_T(y) + b. \mathbf{f} \cdot d_T(b) - y. \mathbf{f} \cdot d_{T^*}(y) - b. \mathbf{f} \cdot d_{T^*}(b) \\ &= x. \mathbf{f} \cdot \frac{d_{T^*}(a)}{d_{T^*}(a)} + a. \mathbf{f} \cdot \frac{d_{T^*}(x)}{d_{T^*}(x)} - x. \mathbf{f} \cdot d_{T^*}(x) - a. \mathbf{f} \cdot d_{T^*}(a) \\ &+ y. \mathbf{f} \cdot \frac{d_{T^*}(b)}{d_{T^*}(b)} + b. \mathbf{f} \cdot \frac{d_{T^*}(y)}{d_{T^*}(y)} - y. \mathbf{f} \cdot d_{T^*}(y) - b. \mathbf{f} \cdot d_{T^*}(b) \\ &= (x. \mathbf{f} - a. \mathbf{f}) (d_{T^*}(a) - d_{T^*}(x)) + (y. \mathbf{f} - b. \mathbf{f}) (d_{T^*}(b) - d_{T^*}(y)) \\ &\leq 0 \end{split}$$

next we show HOFFMAN(C) is correct

Correctness

- \triangleright merge x and y into a new character z
- \triangleright let $C_1 = C x y + z$ and T_1 be optimal for C_1
- \triangleright to show: T obtained from T_1 by splitting z is optimal for C
- \triangleright by definition, $B(T) = B(T_1) + x.f + y.f$





assume
$$B(T') < B(T)$$

assume
$$B(T') < B(T)$$

$$B(T') = B(T'') + x.f + y.f$$

$$\triangleright B(T'') < B(T) - x.f - y.f = B(T_1)$$
, a contradiction