

Huffman's Coding

Tao Hou

- Suppose you need to send a large text to another party
- Let's say, the text consists of six characters, each of which has a frequency of appearing in the text:

<i>sym.</i>	<i>freq.</i>
a	45%
b	13%
c	12%
d	16%
e	9%
f	5%

- We would like to encode the text as a sequence of bits while ensuring that:
 1. we should be able to decode text *without ambiguity*, with an encoding table given
 2. the length of the encoded binary string is *minimal*, or in another word, the average no. of bits for encoding each character is the minimal

Fixed Length Encoding

- A first attempt is to assign each character the same code length
- We have 6 characters, so we need a code of length 3 (bits)
- Such an encoding is referred to as a *fixed-length* encoding

<i>sym.</i>	<i>freq.</i>	<i>code</i>
a	45%	000
b	13%	001
c	12%	010
d	16%	011
e	9%	100
f	5%	101

- Average number of bits to encode a character: 3

Variable-length Encoding

- Can we do better?
- If we can, then we need to forsake the use of fixed-length encoding, and allow character codes to have different lengths
- Such a code is called a *variable-length* code
- Intuitively, we would like to assign characters with high frequencies short codes
- Let's try that

Variable-length Encoding

<i>sym.</i>	<i>freq.</i>	<i>code</i>
a	45%	0
b	13%	1
c	12%	00
d	16%	01
e	9%	10
f	5%	11

- Average no. of bits:

$$.45 \times 1 + .13 \times 1 + .12 \times 2 + .16 \times 2 + .9 \times 2 + .5 \times 2 = 1.42$$

Variable-length Encoding

<i>sym.</i>	<i>freq.</i>	<i>code</i>
a	45%	0
b	13%	1
c	12%	00
d	16%	01
e	9%	10
f	5%	11

- Average no. of bits:

$$.45 \times 1 + .13 \times 1 + .12 \times 2 + .16 \times 2 + .9 \times 2 + .5 \times 2 = 1.42$$

- There is a problem! The above encoding is ambiguous: The text can be decoded in multiple ways

- One way for ensuring decodability in a variable-length code is to require the code to be a *prefix code*: No code is a prefix of another

<i>sym.</i>	<i>freq.</i>	<i>code</i>
a	45%	000
b	13%	001
c	12%	010
d	16%	011
e	9%	10
f	5%	11

Why Prefix Codes Work?

Consider an algorithm for decoding the bit sequence $b[1, \dots, n]$:

1. Maintain a variable i which is the index of the 'next' bit you are scanning at a certain step, i.e., we have scanned all bits $b[1, \dots, i - 1]$ and decoded the texts in them
2. Initially, $i = 0$.
3. Do the following until $i > n$:
 - 3.1 Find an $\ell \geq i$ such that $b[i, \dots, \ell]$ matches a code in the table, and output the corresponding character
 - 3.2 Let $i = \ell + 1$

Why Prefix Codes Work?

Consider an algorithm for decoding the bit sequence $b[1, \dots, n]$:

1. Maintain a variable i which is the index of the ‘next’ bit you are scanning at a certain step, i.e., we have scanned all bits $b[1, \dots, i - 1]$ and decoded the texts in them
2. Initially, $i = 0$.
3. Do the following until $i > n$:
 - 3.1 Find an $\ell \geq i$ such that $b[i, \dots, \ell]$ matches a code in the table, and output the corresponding character
 - 3.2 Let $i = \ell + 1$

We can use proof by contradiction to show that the above process has no ambiguity (i.e., each step finds a unique ℓ) if the bits are encoded with a prefix code

Problem Definition

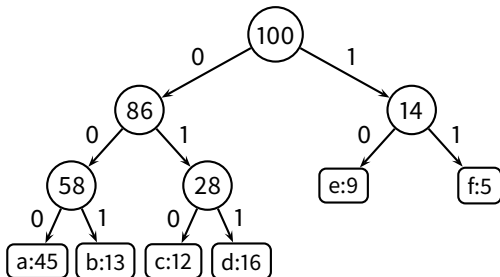
Given a set of characters C and a frequency function $f : C \rightarrow \mathbb{R}$, compute a prefix-code that minimizes

$$\sum_{c \in C} f(c) \times \text{code-length}(c)$$

Prefix Codes & Binary Trees

- Any prefix code can be represented as a binary tree T whose leaves represent the characters
- An edge between a tree node and its left (resp. right) child is labeled 0 (resp. 1)
- the code of a character is the concatenation of the labels on the path from the root to the leaf corresponding to the character

<i>sym.</i>	<i>freq.</i>	<i>code</i>
a	45%	000
b	13%	001
c	12%	010
d	16%	011
e	9%	10
f	5%	11



More on Prefix Codes & Binary Trees

- The algorithm to build such a binary tree: For each character c , we can build a path starting with the root, ending with c , with the edges going to left or right based on the 0/1 in the code of c
- We then combine all such paths to get the binary tree
- Since no code of a character is the prefix of another, we have that no node of a character in the tree is an ancestor of another \Rightarrow all characters correspond to leaves
- So, there is a **one-to-one** correspondence from all such prefix codes to all the binary trees

Problem Reformulation

Given a set of characters C and a frequency function $f : C \rightarrow \mathbb{N}$, construct a binary tree T whose leaves are the characters that minimizes

$$\text{Cost}(T) = \sum_{c \in \text{leaves}(T)} f(c) \times \text{depth}(c)$$

Huffman Algorithm: Ideas

- We work on a bunch of binary trees, whose leaves correspond to the characters
- Each tree has a *frequency*, which is the sum of the frequency of all characters its leaves correspond to

Huffman Algorithm: Ideas

- We work on a bunch of binary trees, whose leaves correspond to the characters
- Each tree has a *frequency*, which is the sum of the frequency of all characters its leaves correspond to
- Initially, each character corresponds to a tree with a single node (frequency is obvious)

Huffman Algorithm: Ideas

- We work on a bunch of binary trees, whose leaves correspond to the characters
- Each tree has a **frequency**, which is the sum of the frequency of all characters its leaves correspond to
- Initially, each character corresponds to a tree with a single node (frequency is obvious)
- Then, in each iteration, we extract two trees with the minimum frequency, and merge the two tree into one by letting them sharing a new root (the frequency will be summed):
This is the **greedy choice**
- We do this until we are left with only a single tree

HUFFMAN(C, f)

```
1  Q = empty heap
2  for each  $c \in C$ 
3      create an isolated tree node  $z$ 
4       $z.char = c$ 
5       $z.freq = f(c)$ 
6      HEAP-INSERT( $Q, z, z.freq$ )
7  for  $i = 1$  to  $n - 1$ 
8      create a new tree node  $z$ 
9       $z.left = \mathbf{EXTRACT-MIN}(Q)$ 
10      $z.right = \mathbf{EXTRACT-MIN}(Q)$ 
11      $z.freq = z.left.freq + z.right.freq$ 
12     HEAP-INSERT( $Q, z, z.freq$ )
13  return EXTRACT-MIN( $Q$ )
```

HUFFMAN(C, f)

```
1  Q = empty heap
2  for each  $c \in C$ 
3      create an isolated tree node  $z$ 
4       $z.char = c$ 
5       $z.freq = f(c)$ 
6      HEAP-INSERT( $Q, z, z.freq$ )
7  for  $i = 1$  to  $n - 1$ 
8      create a new tree node  $z$ 
9       $z.left = \mathbf{EXTRACT-MIN}(Q)$ 
10      $z.right = \mathbf{EXTRACT-MIN}(Q)$ 
11      $z.freq = z.left.freq + z.right.freq$ 
12     HEAP-INSERT( $Q, z, z.freq$ )
13 return EXTRACT-MIN( $Q$ )
```

- Time complexity:

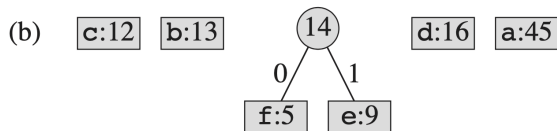
```
HUFFMAN( $C, f$ )
1   $Q =$  empty heap
2  for each  $c \in C$ 
3      create an isolated tree node  $z$ 
4       $z.char = c$ 
5       $z.freq = f(c)$ 
6      HEAP-INSERT( $Q, z, z.freq$ )
7  for  $i = 1$  to  $n - 1$ 
8      create a new tree node  $z$ 
9       $z.left =$  EXTRACT-MIN( $Q$ )
10      $z.right =$  EXTRACT-MIN( $Q$ )
11      $z.freq = z.left.freq + z.right.freq$ 
12     HEAP-INSERT( $Q, z, z.freq$ )
13  return EXTRACT-MIN( $Q$ )
```

- Time complexity: $O(n \log n)$

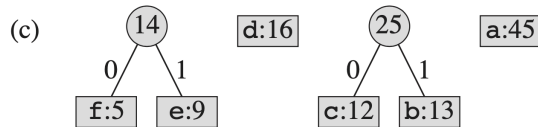
(a)

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

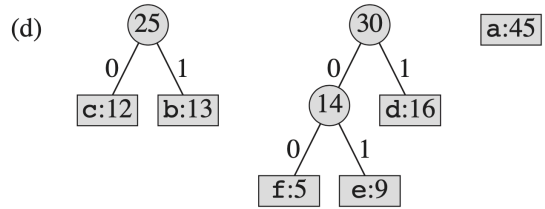
(Example from CLRS)



(Example from CLRS)



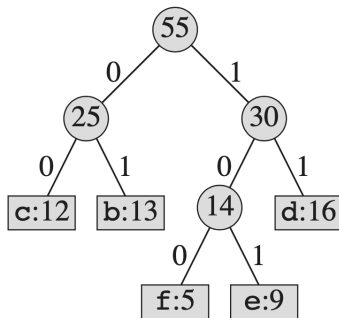
(Example from CLRS)



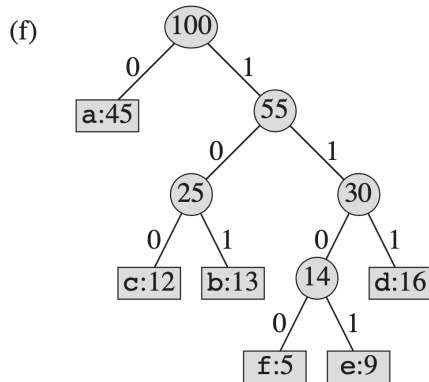
(Example from CLRS)

(e)

a:45



(Example from CLRS)



(Example from CLRS)

Correctness proof of Huffman's algorithm

Lemma 1

Assumption:

- Let C be an alphabet where each character $c \in C$ has frequency $f(c)$.
- Let x and y be two characters in C with minimum frequency.
- Construct another alphabet C' from C by removing x and y and adding a new character z .
 - ▶ Define f for C' as for C , except that $f(z) = f(x) + f(y)$.
- Let T' be any binary tree representing an optimal prefix code for the alphabet C' .

Conclusion:

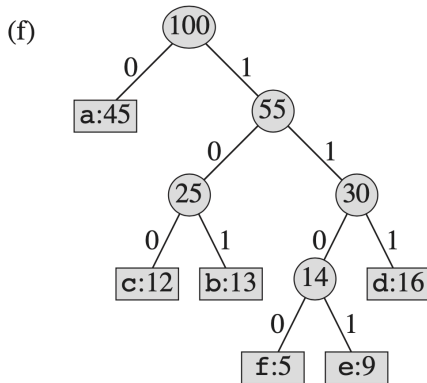
- Then the tree T , obtained from T' by replacing the leaf node for z with an internal node having x and y as children, represents an optimal prefix code for the alphabet C .

Implication of the Lemma

- For the tree T built by Huffman's algorithm, starting from the simplest form containing a single node, we can iteratively expand the tree by replacing a leaf with an internal node plus two leaves, and eventually arrive at T s.t. each tree we have along the expansion is an optimal tree for codes we have.

Implication of the Lemma

- For the tree T built by Huffman's algorithm, starting from the simplest form containing a single node, we can iteratively expand the tree by replacing a leaf with an internal node plus two leaves, and eventually arrive at T s.t. each tree we have along the expansion is an optimal tree for codes we have.
- Let's work on the example using board:



Proving Lemma 1 needs two propositions.

Proposition 1

- Let C be an alphabet where each character $c \in C$ has frequency $f(c)$.
- Let x and y be two characters in C with minimum frequency.
- Then there exists an optimal binary tree for C where x and y are siblings (children of the same parent).

Proposition 1

- Let C be an alphabet where each character $c \in C$ has frequency $f(c)$.
- Let x and y be two characters in C with minimum frequency.
- Then there exists an optimal binary tree for C where x and y are siblings (children of the same parent).

Idea of the proof:

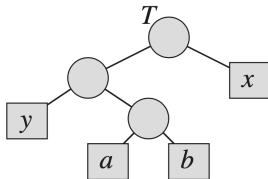
- Take an optimal tree T for C and modify it to make another optimal tree T'' where x and y are sibling leaves of maximum depth.

Proof of Proposition 1

- Let a and b be two characters that are sibling leaves of maximum depth in T .
- WLOG, we can assume $f(a) \leq f(b)$ and $f(x) \leq f(y)$.
- Since $f(x)$ and $f(y)$ are the two min frequencies, we have $f(x) \leq f(a)$ and $f(y) \leq f(b)$.

Proof of Proposition 1

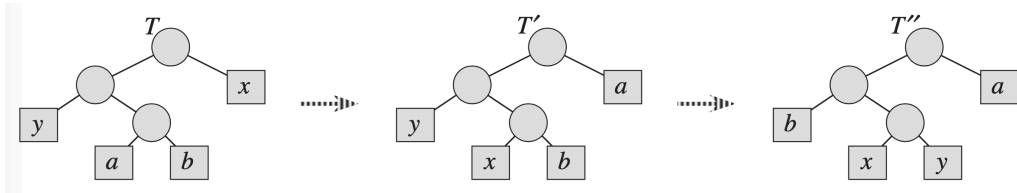
- Let a and b be two characters that are sibling leaves of maximum depth in T .
- WLOG, we can assume $f(a) \leq f(b)$ and $f(x) \leq f(y)$.
- Since $f(x)$ and $f(y)$ are the two min frequencies, we have $f(x) \leq f(a)$ and $f(y) \leq f(b)$.
- To make the arguments easier, assume $x, y \neq a, b$, e.g.:



- ▶ The case of (one of) x, y could be equal to (one of) a, b is trickier and we omit (focus on the big picture)

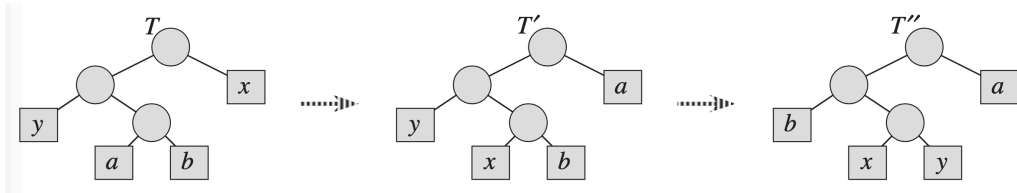
Proof of Proposition 1

- We switch x with a , y with b , to produce the tree T'' we want, where the avg code length cannot increase, so that T'' is also an optimal binary tree:



Proof of Proposition 1

- We switch x with a , y with b , to produce the tree T'' we want, where the avg code length cannot increase, so that T'' is also an optimal binary tree:



■

$$\text{Cost}(T) - \text{Cost}(T') = \sum_{c \in C} f(c)d_T(c) - \sum_{c \in C} f(c)d_{T'}(c) \quad (1)$$

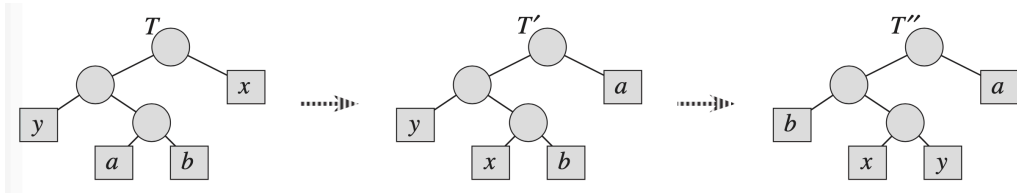
$$= f(x)d_T(x) + f(a)d_T(a) - f(x)d_{T'}(x) - f(a)d_{T'}(a) \quad (2)$$

$$= f(x)d_T(x) + f(a)d_T(a) - f(x)d_T(a) - f(a)d_T(x) \quad (3)$$

$$= (f(a) - f(x))(d_T(a) - d_T(x)) \geq 0 \quad (4)$$

Proof of Proposition 1

- We switch x with a , y with b , to produce the tree T'' we want, where the avg code length cannot increase, so that T'' is also an optimal binary tree:



■

$$\text{Cost}(T) - \text{Cost}(T') = \sum_{c \in C} f(c)d_T(c) - \sum_{c \in C} f(c)d_{T'}(c) \quad (1)$$

$$= f(x)d_T(x) + f(a)d_T(a) - f(x)d_{T'}(x) - f(a)d_{T'}(a) \quad (2)$$

$$= f(x)d_T(x) + f(a)d_T(a) - f(x)d_T(a) - f(a)d_T(x) \quad (3)$$

$$= (f(a) - f(x))(d_T(a) - d_T(x)) \geq 0 \quad (4)$$

- $\text{Cost}(T) \geq \text{Cost}(T') \geq \text{Cost}(T'')$

Proposition 2

- Let T' be a binary tree (corresponding to a prefix code) with a leaf z
- Let tree T be obtained from T' by replacing z with an internal node having x and y as children, s.t. $f(z) = f(x) + f(y)$.
- Then, $Cost(T') = Cost(T) - f(x) - f(y)$

Proposition 2

- Let T' be a binary tree (corresponding to a prefix code) with a leaf z
- Let tree T be obtained from T' by replacing z with an internal node having x and y as children, s.t. $f(z) = f(x) + f(y)$.
- Then, $Cost(T') = Cost(T) - f(x) - f(y)$

Proof:

■

$$Cost(T) - Cost(T') = \sum_{c \in C} f(c)d_T(c) - \sum_{c \in C'} f(c)d_{T'}(c) \quad (5)$$

$$= f(x)d_T(x) + f(y)d_T(y) - f(z)d_{T'}(z) \quad (6)$$

Proposition 2

- Let T' be a binary tree (corresponding to a prefix code) with a leaf z
- Let tree T be obtained from T' by replacing z with an internal node having x and y as children, s.t. $f(z) = f(x) + f(y)$.
- Then, $Cost(T') = Cost(T) - f(x) - f(y)$

Proof:



$$Cost(T) - Cost(T') = \sum_{c \in C} f(c)d_T(c) - \sum_{c \in C'} f(c)d_{T'}(c) \quad (5)$$

$$= f(x)d_T(x) + f(y)d_T(y) - f(z)d_{T'}(z) \quad (6)$$

- Since $d_T(x) = d_T(y) = d_{T'}(z) + 1$ and $f(z) = f(x) + f(y)$, We have:

$$Cost(T) - Cost(T') = (f(x) + f(y))(d_{T'}(z) + 1) - (f(x) + f(y))d_{T'}(z) \quad (7)$$

$$= f(x) + f(y) \quad (8)$$

- Use proof by contradiction. Suppose that T is not an optimal binary tree for C .
- Then there exists an optimal tree T'' for C such that $Cost(T'') < Cost(T)$.

- Use proof by contradiction. Suppose that T is not an optimal binary tree for C .
- Then there exists an optimal tree T'' for C such that $Cost(T'') < Cost(T)$.
- Since x and y have minimum frequency in C , we can also let T'' be an optimal tree where x and y are siblings (by Proposition 1).

- Use proof by contradiction. Suppose that T is not an optimal binary tree for C .
- Then there exists an optimal tree T'' for C such that $Cost(T'') < Cost(T)$.
- Since x and y have minimum frequency in C , we can also let T'' be an optimal tree where x and y are siblings (by Proposition 1).
- Let T''' be derived from T'' with the common parent of x and y replaced by a leaf z with $f(z) = f(x) + f(y)$.
 - ▶ T''' is a binary tree for C'

- Use proof by contradiction. Suppose that T is not an optimal binary tree for C .
- Then there exists an optimal tree T'' for C such that $Cost(T'') < Cost(T)$.
- Since x and y have minimum frequency in C , we can also let T'' be an optimal tree where x and y are siblings (by Proposition 1).
- Let T''' be derived from T'' with the common parent of x and y replaced by a leaf z with $f(z) = f(x) + f(y)$.
 - ▶ T''' is a binary tree for C'

$$Cost(T''') \tag{9}$$

$$= Cost(T'') - f(x) - f(y) \tag{10}$$

- Use proof by contradiction. Suppose that T is not an optimal binary tree for C .
- Then there exists an optimal tree T'' for C such that $Cost(T'') < Cost(T)$.
- Since x and y have minimum frequency in C , we can also let T'' be an optimal tree where x and y are siblings (by Proposition 1).
- Let T''' be derived from T'' with the common parent of x and y replaced by a leaf z with $f(z) = f(x) + f(y)$.
 - ▶ T''' is a binary tree for C'
-

$$Cost(T''') \tag{9}$$

$$= Cost(T'') - f(x) - f(y) \tag{10}$$

$$< Cost(T) - f(x) - f(y) \tag{11}$$

- Use proof by contradiction. Suppose that T is not an optimal binary tree for C .
- Then there exists an optimal tree T'' for C such that $Cost(T'') < Cost(T)$.
- Since x and y have minimum frequency in C , we can also let T'' be an optimal tree where x and y are siblings (by Proposition 1).
- Let T''' be derived from T'' with the common parent of x and y replaced by a leaf z with $f(z) = f(x) + f(y)$.
 - ▶ T''' is a binary tree for C'

$$Cost(T''') \tag{9}$$

$$= Cost(T'') - f(x) - f(y) \tag{10}$$

$$< Cost(T) - f(x) - f(y) \tag{11}$$

$$= Cost(T') \tag{12}$$

where (9)-(10), (11)-(12) are by Proposition 2.

- This contradicts the fact that T' is an optimal binary tree for C'

- Certain figures in the slides are taken from [CLRS]