

Greedy algorithms

Huffman Encoding (Textbook Section 5.2)

Huffman Encoding

An encoding scheme used in, e.g., MP3 encoding

Huffman Encoding

An encoding scheme used in, e.g., MP3 encoding

Data: a string S of symbols over an alphabet Γ

Huffman Encoding

An encoding scheme used in, e.g., MP3 encoding

Data: a string S of symbols over an alphabet Γ

Goal: find a binary encoding e of Γ resulting in minimum encoded length of S

Huffman Encoding

An encoding scheme used in, e.g., MP3 encoding

Data: a string S of symbols over an alphabet Γ

Goal: find a binary encoding e of Γ resulting in minimum encoded length of S

Denote the encoded string by S_e

Huffman Encoding

$$S = abaa$$

An encoding scheme used in, e.g., MP3 encoding

Data: a string S of symbols over an alphabet Γ

Goal: find a binary encoding e of Γ resulting in minimum encoded length of S

Denote the encoded string by S_e

$$S_e = \underbrace{01100001}_a \underbrace{01100010}_b \underbrace{01100001}_a$$

Example: ASCII encoding

a	01100001
b	01100010
\vdots	\vdots

e''

Different encodings

Consider $\Gamma = \{a, b, c\}$

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$a \rightarrow 00$

$e_1 : b \rightarrow 01$

$c \rightarrow 10$

$$|S_{e_1}| = 45 \cdot 2 + 16 \cdot 2 + 2 \cdot 2 = 126$$

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$$a \rightarrow 00$$

$$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$$

$$c \rightarrow 10$$

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$$a \rightarrow 00$$

$$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$$

$$c \rightarrow 10$$

- Variable-length encoding

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$a \rightarrow 00$

$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$

$c \rightarrow 10$

- Variable-length encoding

$a \rightarrow 0$

$e_2 : b \rightarrow 10$

$c \rightarrow 11$

$$|S_{e_2}| = 45 \cdot 1 + 16 \cdot 2 + 2 \cdot 2 =$$

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$$a \rightarrow 00$$

$$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$$

$$c \rightarrow 10$$

- Variable-length encoding

$$a \rightarrow 0$$

$$e_2 : b \rightarrow 10 \quad |S_{e_2}| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

$$c \rightarrow 11$$

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$a \rightarrow 00$

$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$

$c \rightarrow 10$

- Variable-length encoding

$a \rightarrow 0$

$e_2 : b \rightarrow 10 \quad |S_{e_2}| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$

$c \rightarrow 11$

- Be careful!

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$$a \rightarrow 00$$

$$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$$

$$c \rightarrow 10$$

- Variable-length encoding

$$a \rightarrow 0$$

$$e_2 : b \rightarrow 10 \quad |S_{e_2}| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

$$c \rightarrow 11$$

$$a \rightarrow 0$$

- Be careful! $e_2 : b \rightarrow 1$

$$c \rightarrow 01$$

Different encodings

Consider $\Gamma = \{a, b, c\}$

Stats on S : a appears 45 times, b 16 times, and c twice

- Fixed-length encoding

$$a \rightarrow 00$$

$$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$$

$$c \rightarrow 10$$

- Variable-length encoding

$$a \rightarrow 0$$

$$e_2 : b \rightarrow 10 \quad |S_{e_2}| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

$$c \rightarrow 11$$

$$a \rightarrow 0$$

- Be careful! $e_2 : b \rightarrow 1$ Decoding will lead to ambiguity

$$c \rightarrow 01$$

Prefix-free encoding

$$a \rightarrow 0$$

Consider the bad encoding e_2 :

$$b \rightarrow 1$$

$$c \rightarrow 01$$

Prefix-free encoding

Consider the bad encoding e_2 :

$a \rightarrow 0$	
$b \rightarrow 1$	
$c \rightarrow 01$	

How to decode $\overbrace{01}^c \overbrace{01}^a \overbrace{10}^b$?

Prefix-free encoding

$$a \rightarrow 0$$

Consider the bad encoding e_2 : $b \rightarrow 1$ How to decode 010110?

$$c \rightarrow 01$$

ababba?, *ccba?*, *abcba?*, or ...?

Prefix-free encoding

$$a \rightarrow 0$$

Consider the bad encoding $e_2 : b \rightarrow 1$ How to decode 010110?

$$c \rightarrow 01$$

ababba?, *ccba?*, *abcba?*, or ...?

To avoid ambiguity, we need the encoding to be **prefix-free**

Prefix-free encoding

$$a \rightarrow 0$$

Consider the bad encoding $e_2 : b \rightarrow 1$ How to decode 010110?

$$c \rightarrow 01$$

ababba?, *ccba?*, *abcba?*, or ...?

To avoid ambiguity, we need the encoding to be **prefix-free**

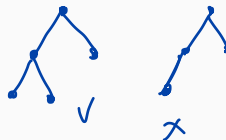
Definition

An encoding is **prefix-free** if no codeword is a prefix of any other codewords

Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

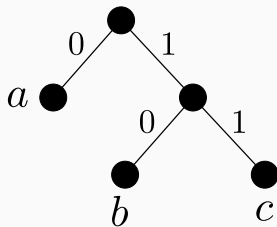


Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding



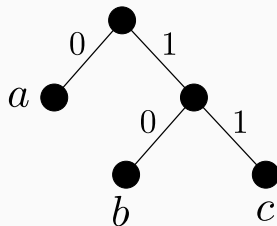
Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

- leaves are corresponding to symbols in Γ



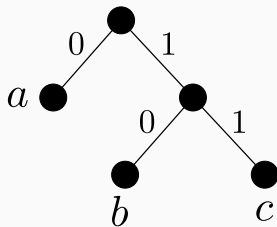
Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

- leaves are corresponding to symbols in Γ
- label edge to the left child with 0



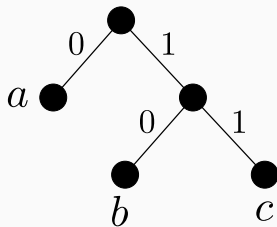
Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

- leaves are corresponding to symbols in Γ
- label edge to the left child with 0
- label edge to the right child with 1



Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

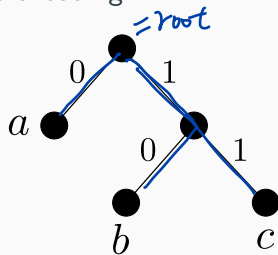
- leaves are corresponding to symbols in Γ
- label edge to the left child with 0
- label edge to the right child with 1

To obtain the encoding, read edge labels from root to a symbol

a: 0

b: 10

c: 11



Tree representation of a prefix-free encoding (I)

Definition

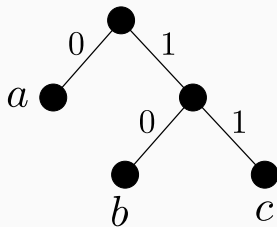
A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

- leaves are corresponding to symbols in Γ
- label edge to the left child with 0
- label edge to the right child with 1

To obtain the encoding, read edge labels from root to a symbol

$a \rightarrow 0$, $b \rightarrow 10$, $c \rightarrow 11$,



Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

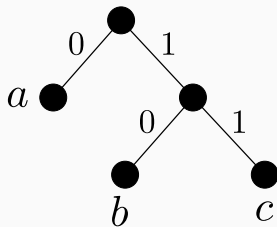
We use a full binary tree to represent a prefix-free encoding

- leaves are corresponding to symbols in Γ
- label edge to the left child with 0
- label edge to the right child with 1

To obtain the encoding, read edge labels from root to a symbol

$a \rightarrow 0$, $b \rightarrow 10$, $c \rightarrow 11$,

Depth of a leaf \equiv length of its codeword



Tree representation of a prefix-free encoding (I)

Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

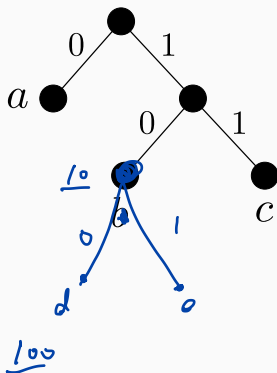
- leaves are corresponding to symbols in Γ
- label edge to the left child with 0
- label edge to the right child with 1

To obtain the encoding, read edge labels from root to a symbol

$a \rightarrow 0$, $b \rightarrow 10$, $c \rightarrow 11$,

Depth of a leaf \equiv length of its codeword

It guarantees to be prefix-free



Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

Tree representation of a prefix-free encoding (II)

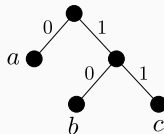
Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$



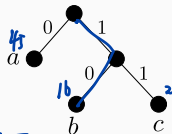
Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$

the number of bits to represent v



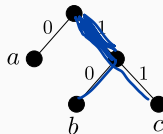
$$|S_e| = 45 \cdot 1 + 16 \cdot 2 + 2 \cdot 2 = 81$$

Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



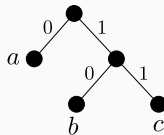
$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

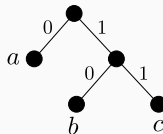
A useful re-write: label internal nodes with counts of descendants

Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

A useful re-write: label internal nodes with counts of descendants

For all non-root node v , define

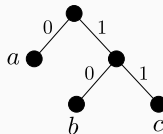
$\text{cost}(v) := \text{sum of leaf node counts descending from } v$

Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$

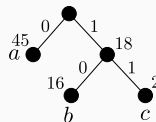


$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

A useful re-write: label internal nodes with counts of descendants

For all non-root node v , define

$\text{cost}(v) := \text{sum of leaf node counts descending from } v$

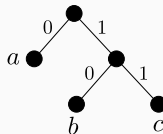


Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

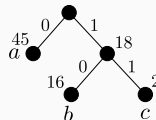
A useful re-write: label internal nodes with counts of descendants

For all non-root node v , define

$\text{cost}(v) :=$ sum of leaf node counts descending from v

$$|S_e| = \sum_{v \in T - \{\text{root}\}} \text{cost}(v)$$

all nodes except for the root.



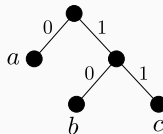
$$|S_e| = 45 + 16 + 2 + 18$$

Tree representation of a prefix-free encoding (II)

Let e be an encoding represented by a tree

For string S , let f_v be the symbol count in S for each $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



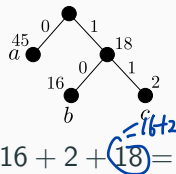
$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

A useful re-write: label internal nodes with counts of descendants

For all non-root node v , define

$\text{cost}(v) :=$ sum of leaf node counts descending from v

$$|S_e| = \sum_{v \in T - \{\text{root}\}} \text{cost}(v)$$



$$|S_e| = 45 + 16 + 2 + 18 = 81$$

Constructing the prefix-free encoding tree

Idea: put more frequent symbols at smaller depth

Constructing the prefix-free encoding tree

Idea: put more frequent symbols at smaller depth

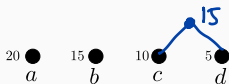
Greedy approach: continually merge least frequent symbols/nodes until you have a full binary tree encoding all symbols

Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$

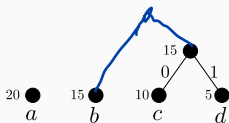
Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



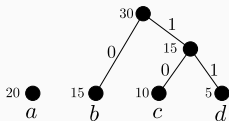
Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



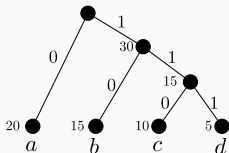
Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



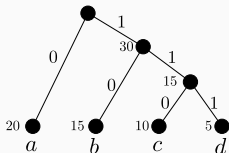
Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



$a \rightarrow 0$

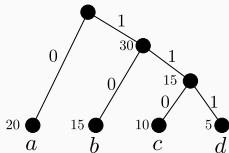
$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



$a \rightarrow 0$

$b \rightarrow 10$

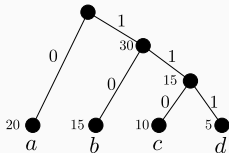
$c \rightarrow 110$

$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$

Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

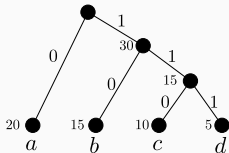
$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



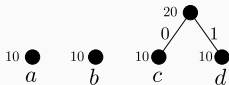
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

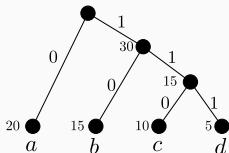
$c \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



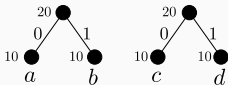
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

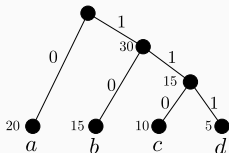
$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



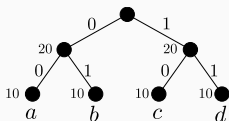
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

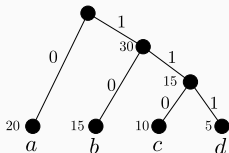
$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



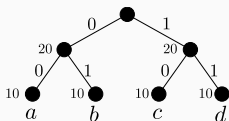
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



$a \rightarrow 00$

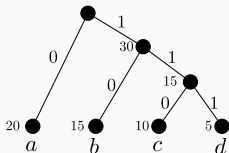
$b \rightarrow 01$

$c \rightarrow 10$

$d \rightarrow 11$

Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



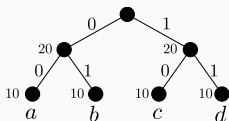
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$c \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



$a \rightarrow 00$

$b \rightarrow 01$

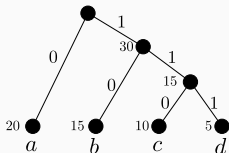
$c \rightarrow 10$

$c \rightarrow 11$

- $a : 10, b : 10, c : 5, d : 5$

Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



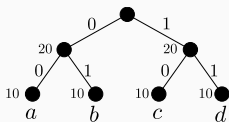
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$c \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



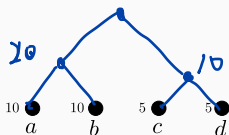
$a \rightarrow 00$

$b \rightarrow 01$

$c \rightarrow 10$

$c \rightarrow 11$

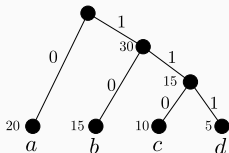
- $a : 10, b : 10, c : 5, d : 5$



$$|S_e| = 10 + 10 + 5 + 5 + 20 + 10 = 60$$

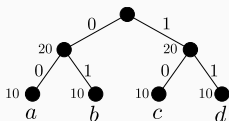
Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



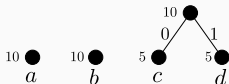
$a \rightarrow 0$
 $b \rightarrow 10$
 $c \rightarrow 110$
 $d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



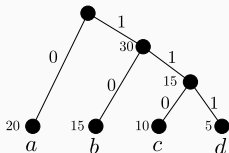
$a \rightarrow 00$
 $b \rightarrow 01$
 $c \rightarrow 10$
 $d \rightarrow 11$

- $a : 10, b : 10, c : 5, d : 5$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



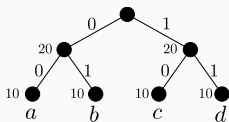
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



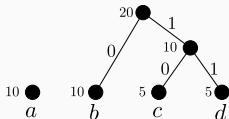
$a \rightarrow 00$

$b \rightarrow 01$

$c \rightarrow 10$

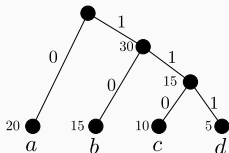
$d \rightarrow 11$

- $a : 10, b : 10, c : 5, d : 5$



Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



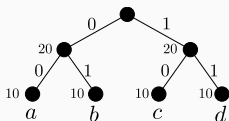
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



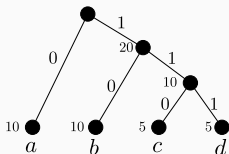
$a \rightarrow 00$

$b \rightarrow 01$

$c \rightarrow 10$

$d \rightarrow 11$

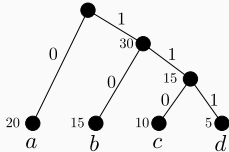
- $a : 10, b : 10, c : 5, d : 5$



$$|S_e| = 10 + 10 + 5 + 5 + 10 + 20 = 60$$

Constructing the prefix-free encoding tree – examples

- $a : 20, b : 15, c : 10, d : 5$



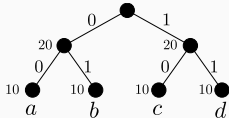
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



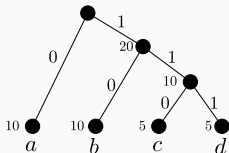
$a \rightarrow 00$

$b \rightarrow 01$

$c \rightarrow 10$

$d \rightarrow 11$

- $a : 10, b : 10, c : 5, d : 5$



$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

Proof of optimality (I)

Proof sketch

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Suppose we have a tree T with two lowest frequent symbols not as deep as possible.

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Suppose we have a tree T with two lowest frequent symbols not as deep as possible. Then at least one has a smaller depth.

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Suppose we have a tree T with two lowest frequent symbols not as deep as possible. Then at least one has a smaller depth. Switch it with one of the deepest nodes that is more frequent.

Proof of optimality (I)

Proof sketch

- Greedy choice property:

Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Suppose we have a tree T with two lowest frequent symbols not as deep as possible. Then at least one has a smaller depth. Switch it with one of the deepest nodes that is more frequent.

This improves the encoding length. Thus T is not optimal



Proof of optimality (I)

Proof sketch

- Greedy choice property:

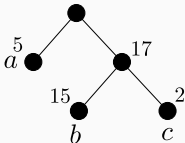
Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Suppose we have a tree T with two lowest frequent symbols not as deep as possible. Then at least one has a smaller depth. Switch it with one of the deepest nodes that is more frequent.

This improves the encoding length. Thus T is not optimal □



Proof of optimality (I)

Proof sketch

- Greedy choice property:

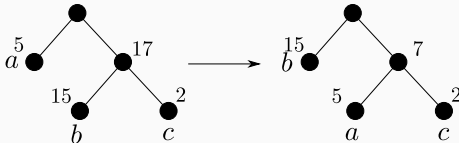
Claim

Every optimal solution has two lowest frequent symbols as leaves connected to an internal node of greatest depth

Proof. (exchange argument).

Suppose we have a tree T with two lowest frequent symbols not as deep as possible. Then at least one has a smaller depth. Switch it with one of the deepest nodes that is more frequent.

This improves the encoding length. Thus T is not optimal □



Proof of optimality (II)

- Optimal substructure

Consider a global optimal solution for (f_1, f_2, \dots, f_n) with f_1, f_2 the least frequencies.

Proof of optimality (II)

- Optimal substructure

Consider a global optimal solution for (f_1, f_2, \dots, f_n) with f_1, f_2 the least frequencies.

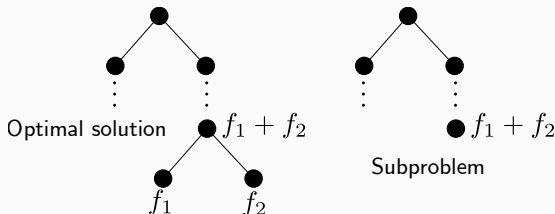
Our subproblem is: find an optimal solution for $(f_1 + f_2, f_3, \dots, f_n)$

Proof of optimality (II)

- Optimal substructure

Consider a global optimal solution for (f_1, f_2, \dots, f_n) with f_1, f_2 the least frequencies.

Our subproblem is: find an optimal solution for $(f_1 + f_2, f_3, \dots, f_n)$

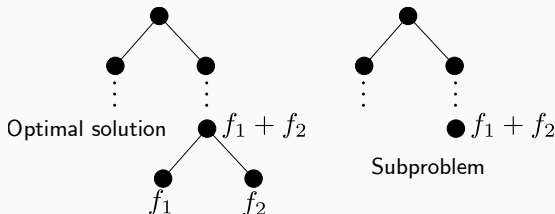


Proof of optimality (II)

- Optimal substructure

Consider a global optimal solution for (f_1, f_2, \dots, f_n) with f_1, f_2 the least frequencies.

Our subproblem is: find an optimal solution for $(f_1 + f_2, f_3, \dots, f_n)$



Thus, the greedy solution will lead to the global optimal solution

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
```

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
```

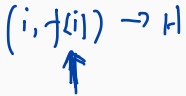
Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
    |
```

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
```

$(i, f[i]) \rightarrow h$



Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
      $\lfloor$ 
```


Pseudocode

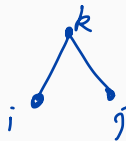
```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
```

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
```

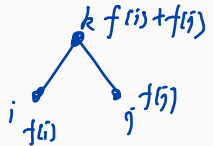
Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
9     Creat a node  $k$  in  $T$  with children  $i$  and  $j$ ;
```



Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
9     Creat a node  $k$  in  $T$  with children  $i$  and  $j$ ;
10     $f[k] := f[i] + f[j]$ ;
```



Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
9     Creat a node  $k$  in  $T$  with children  $i$  and  $j$ ;
10     $f[k] := f[i] + f[j]$ ;
11     $\lfloor$  insert( $H, k$ );
```

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\text{insert}(H, i)$ ;  $\rightarrow O(\log n)$  }  $O(n \log n)$ 
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;  $\rightarrow O(\log n)$ 
8      $j := \text{extract\_min}(H)$ ;  $\rightarrow O(\log n)$  }  $O(n \log n)$ 
9     Create a node  $k$  in  $T$  with children  $i$  and  $j$ ;
10     $f[k] := f[i] + f[j]$ ;
11     $\text{insert}(H, k)$ ;  $\rightarrow O(\log n)$ 
```

Binary heap: insert $O(\log n)$, extract_min: $O(\log n)$

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
9     Creat a node  $k$  in  $T$  with children  $i$  and  $j$ ;
10     $f[k] := f[i] + f[j]$ ;
11     $\lfloor$  insert( $H, k$ );
```

Binary heap: insert $O(\log n)$, extract_min: $O(\log n)$

Lines 4-5: $O(n \log n)$;

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
9     Creat a node  $k$  in  $T$  with children  $i$  and  $j$ ;
10     $f[k] := f[i] + f[j]$ ;
11     $\lfloor$  insert( $H, k$ );
```

Binary heap: insert $O(\log n)$, extract_min: $O(\log n)$

Lines 4-5: $O(n \log n)$; Lines 6-11: $O(n \log n)$

Pseudocode

```
1 def HUFFMAN( $f$ ): //  $f : f[1], \dots, f[n]$ ;  $\Gamma$  has  $n$  symbols
2    $T$ : empty tree;
3    $H$ : priority queue ordered by  $f$ ;
4   for  $i := 1$  in  $n$ :
5      $\lfloor$  insert( $H, i$ );
6   for  $k := n + 1$  in  $2n - 1$ :
7      $i := \text{extract\_min}(H)$ ;
8      $j := \text{extract\_min}(H)$ ;
9     Creat a node  $k$  in  $T$  with children  $i$  and  $j$ ;
10     $f[k] := f[i] + f[j]$ ;
11     $\lfloor$  insert( $H, k$ );
```

Binary heap: insert $O(\log n)$, extract_min: $O(\log n)$

Lines 4-5: $O(n \log n)$; Lines 6-11: $O(n \log n)$

Total cost: $O(n \log n)$

More about the pseudocode

Question: why $2n - 1$ in line 6?

More about the pseudocode

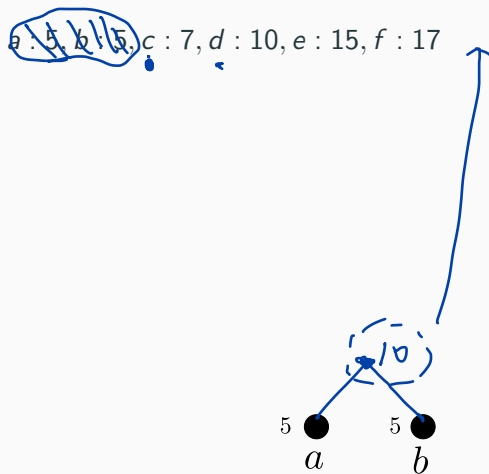
Question: why $2n - 1$ in line 6?

Answer: if a full binary tree has n leaves, then it has $2n - 1$ total nodes

More examples

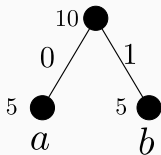
$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$

More examples



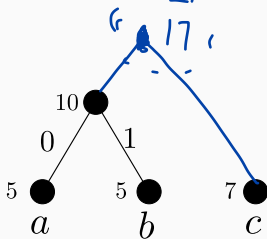
More examples

$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$



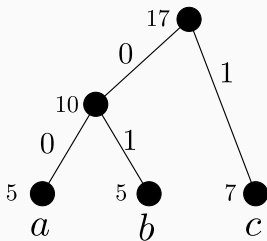
More examples

~~$a : 5, b : 5, c : 7$~~ , $d : 10, e : 15, f : 17$



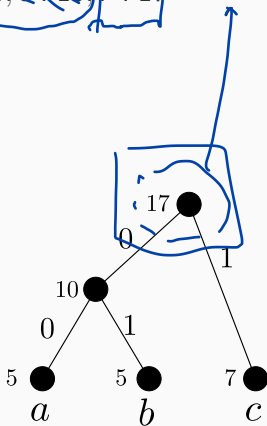
More examples

$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$



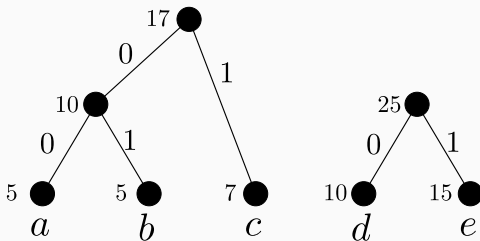
More examples

~~$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$~~



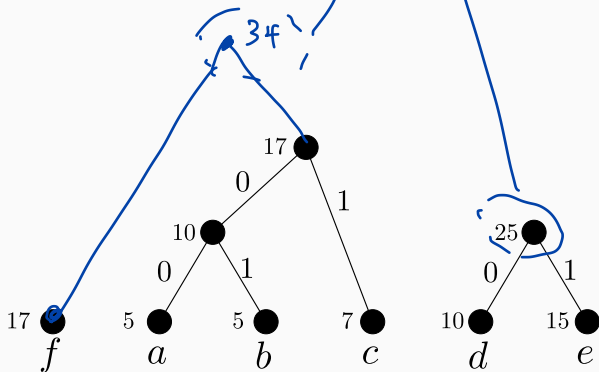
More examples

$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$



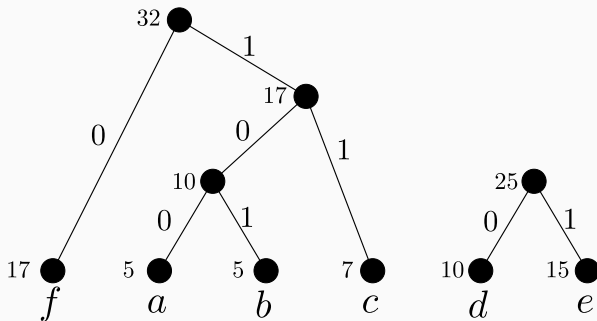
More examples

~~$a: 5, b: 5, c: 7, d: 10, e: 15, f: 17$~~



More examples

$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$



More examples

$a : 5, b : 5, c : 7, d : 10, e : 15, f : 17$

