

LECTURE 10 – GREEDY ALGORITHMS (PART 2)

COMPSCI 308 – DESIGN AND ANALYSIS OF ALGORITHMS

Xing Shi Cai <https://newptcai.gitlab.io>

Duke Kunshan University, Kunshan, China

Jan-Mar 2026

ITA 15.3 Huffman codes

Correctness of Huffman's Algorithm

ASSIGNMENTS¹



Practice makes perfect!

Introduction to Algorithms (ITA) 
Required Readings:

- Section 15.3.
-  Required Exercises:
- Exercises 15.3: 1–8.

¹ ⚪ Assignments will not be collected; however, quiz problems will be selected from them. (This includes both Readings and Exercises.)

ITA 15.3 HUFFMAN CODES

How to TRAVEL IN SPACE

It is extremely difficult to built 🚀 to take humans to traverse the interstellar space.

In the SF trilogy, *Three Bodies*, earth sent a 🚀 with a human 🧠 to intercept 👽 invaders.

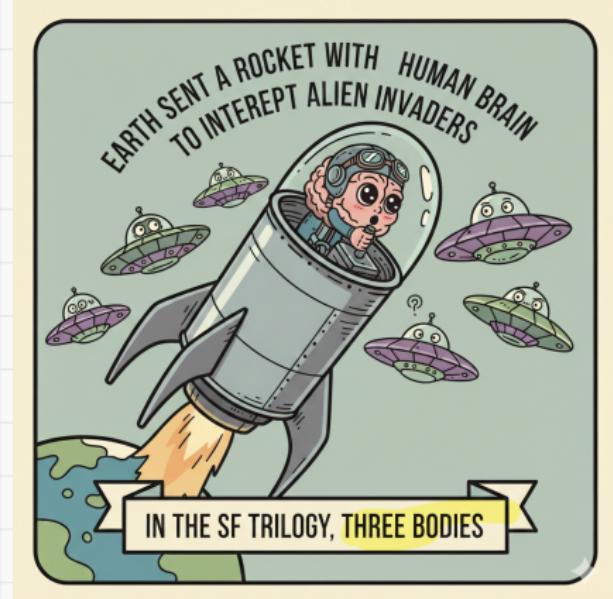


Figure 1: A 🚀 sending a 🧠 to 👽

How to TRAVEL IN SPACE

It is extremely difficult to built 🚀 to take humans to traverse the interstellar space.

In the SF trilogy, *Three Bodies*, earth sent a 🚀 with a human 🧑 to intercept 👽 invaders.

Nowadays they could send digital copies of a human, such as ChatGPT instead instead.

🤔 How can we compress the data so we can send as many digital humans as possible?



Figure 1: Digital humans travelling through space

ENCODING DIGITAL DATA

Consider the text:

✓ {a, b, c, d, e, f}

a deaf dad bade a bee feed a babe



ENCODING DIGITAL DATA

Consider the text:

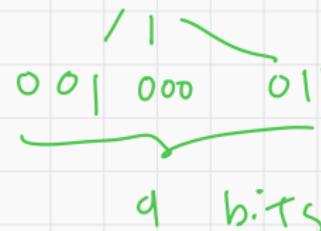
a deaf dad bade a bee feed a babe

To represent the text as a **binary string**, we can use the following **code**:

Letters	a	b	d	e	f
Fixed-Len. Codewords	000	001	011	100	101

We call the binary string representation of a letter its **codeword**.

🎂 How can we encode the word **bad** with the above codeword?



ENCODING DIGITAL DATA

Consider the text:

a deaf dad bade a bee feed a babe

To represent the text as a binary string, we can use the following **code**:

Letters	a	b	d	e	f
Fixed-Len. Codewords	000	001	011	100	101

The example text (with white space removed) is encoded as:

000 011 100 000 101 000 011 011
001 000 011 100 000 001 100 100
101 100 100 011 000 001 000 001 100

The result is a binary string of length 75.

VARIABLE-LENGTH CODEWORD

We also can take the frequency of each letter into account —

more frequent.

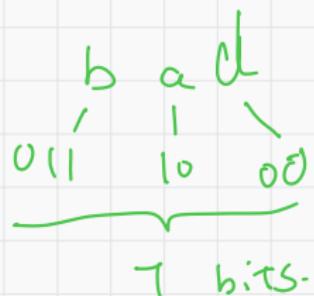
Least frequent.

Character	a	e	d	b	f
Frequency	0.28	0.28	0.2	0.16	0.08
Counts	7	7	5	4	2
Variable-Len. Codewords	10	11	00	011	010

↑ short codeword

long code word

🎂 How can we encode the word *bad* with the above codewords?



VARIABLE-LENGTH CODEWORD

We also can take the frequency of each letter into account –

Character	a	e	d	b	f
Frequency	0.28	0.28	0.2	0.16	0.08
Counts	7	7	5	4	2
Variable-Len. Codewords	10	11	00	011	010

Encoding the text with variable-length codewords, we get:

10 00 11 10 010 10 00 11 10
10 011 10 00 11 10 011 10 00
11 10 011 10 00 11 010 11 10
00 11 10 011 10 00 11 010 11
10 00 11 10 010 10

The result is a binary string of length 56 < 75.

THE COST OF A CODE

Character	a	e	d	b	f
Frequency	0.28	0.28	0.2	0.16	0.08
Counts	7	7	5	4	2
Variable-Len. Codewords	10	11	00	011	010

Prefix code.

With this code, the average length of a code word in the encoded text is

$$0.28 \cdot 2 + 0.28 \cdot 2 + 0.2 \cdot 2 + 0.16 \cdot 3 + 0.08 \cdot 3 = 2.24$$

↑ Average number
of bits

We call this the **cost** of the code.

Our aim is find the **optimal** code, i.e., the one with the minimum cost.

Question: Given binary string.

0|1|0|0

for encoding
a letter.

how to split it into code words.

CAUTION!

Consider the following codewords:

Character	a	e	d	c	b	f
Codewords	10	11	01	001	011	010

🎂 What does the following binary string represent?

01|11|10|001

 ^ ^ ^ ^
 b a d c

 d e c

Ambiguity ↙

PREFIX CODE

To avoid ambiguity, we can will only consider **prefix codes** in which **no** codeword is prefix of another.

... Which ones of the following two are prefix codes?

The first
pure.

Character	a	e	d	c	b	f
Codewords	10	11	00	001	011	010

Character	a	e	d	c	b	f
Codewords	10	11	01	001	011	010

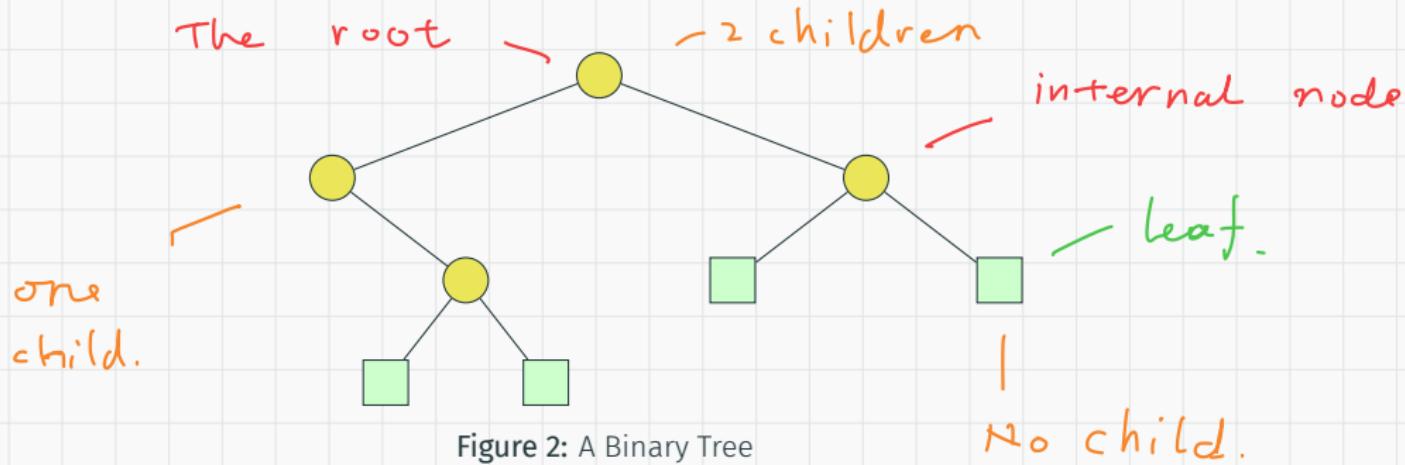
Not prefix code.



BINARY TREE

A **binary tree** is **tree** in which each **node** has at most two children.

A node without any children is called a **leaf**.



BINARY TREE REPRESENTATION

We can represent the following code

Character	a	e	d	b	f
Frequency	0.28	0.28	0.2	0.16	0.08
Counts	7	7	5	4	2
Variable-Len. Codewords	10	11	00	011	010

with a binary tree, with each leaf representing a letter.

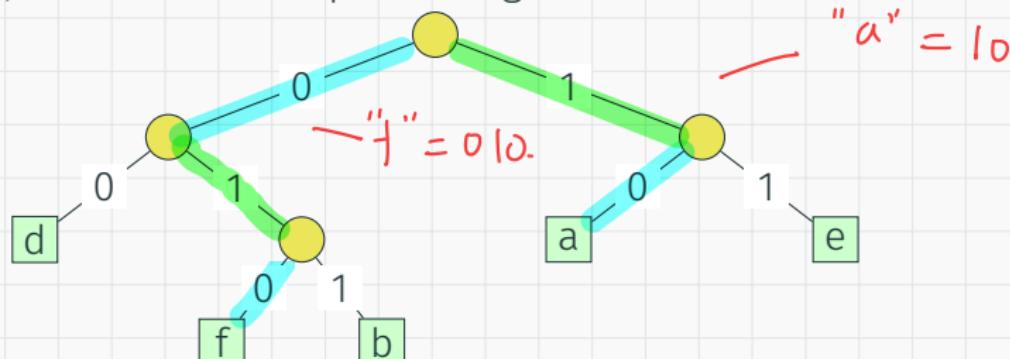
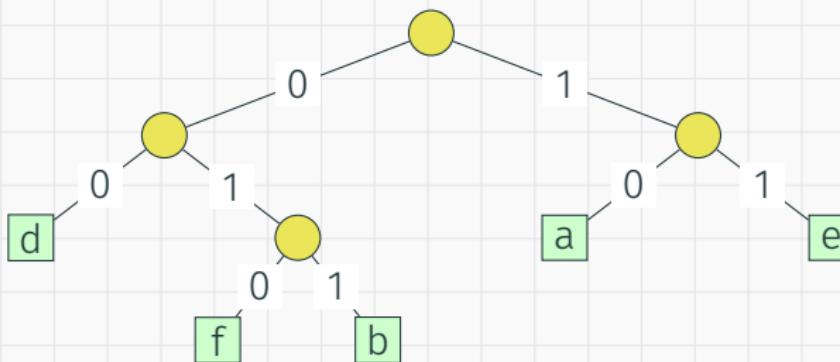


Figure 3: A Binary Tree Representation of the above codewords



BINARY TREE REPRESENTATION OF PREFIX CODES

Do you see why a code with a binary tree representation is a prefix code?



Give such
a tree,
we have
a prefix
code

Figure 4: Example of a full binary tree T representation of a code

The two paths going to two different leaves, must split at some point, so one code word can not be the prefix of the other.



BINARY TREE REPRESENTATION OF PREFIX CODES

Do you see why a code with a binary tree representation is a prefix code?

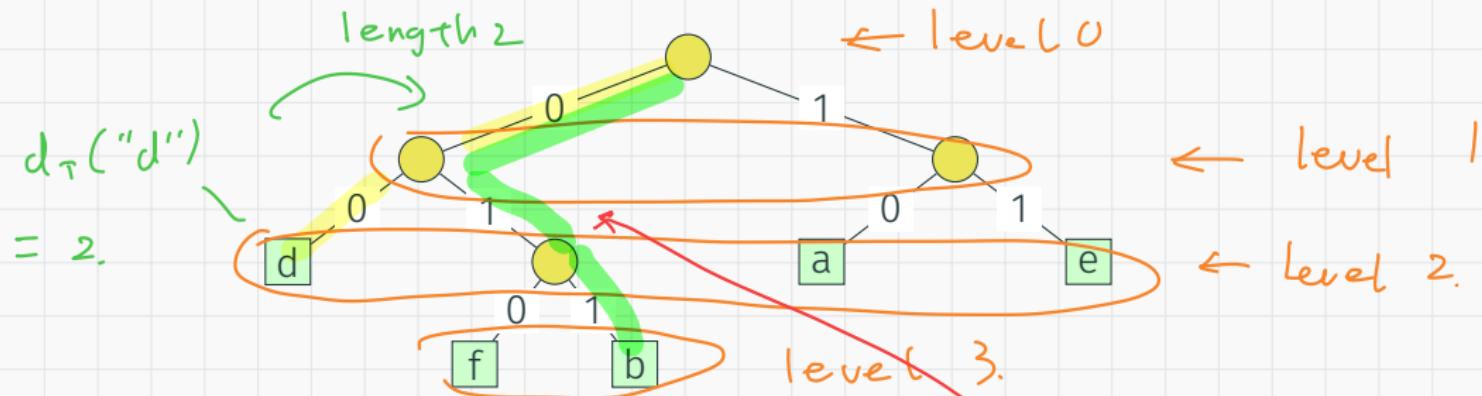


Figure 4: Example of a full binary tree T representation of a code

For a letter ℓ , let $d_T(\ell)$ be the depth of the ℓ representing the c .

What are the $d_T(\cdot)$ in the above tree?

$$d_T(b) = 3.$$

THE COST REVISED

Let C be the set of letters in a text.

For $c \in C$, let $c.freq$ be its frequency.

The cost of a code can be written as

$$\sum_{c \in C} d_T(c) \cdot c.freq,$$



The length of c 's code word.

THE COST REVISED

The **cost** of a code can be written as

$$\sum_{c \in C} d_T(c) \cdot c.freq,$$

$$2 \rightarrow 0.28 + 3 \cdot 0.16$$

+ ...



What is the cost for the following tree/code?

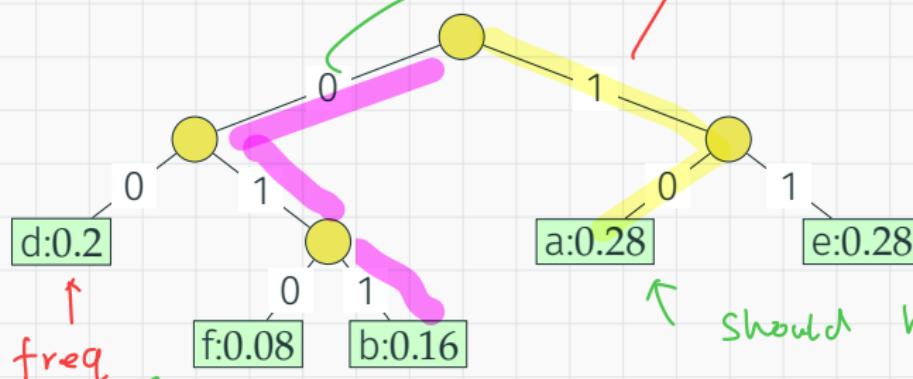


Figure 5: Example of a tree representing a code

should be close to
root.

Be away from the root.

FULL BINARY TREE

A **full binary tree** is a tree in which each node has either two or zero children.

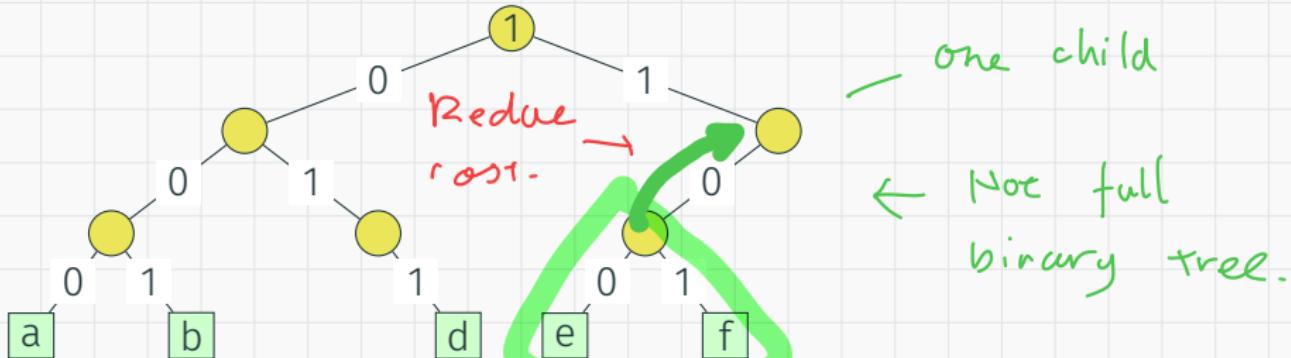
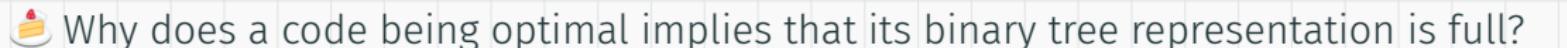


Figure 6: Not a full binary tree, not an optimal code

For a code to be optimal, it is necessary to have full binary tree.

NUMBER OF LEAVES

A full binary tree with n leaves has $n - 1$ internal nodes.

TRY Try to prove this by induction.

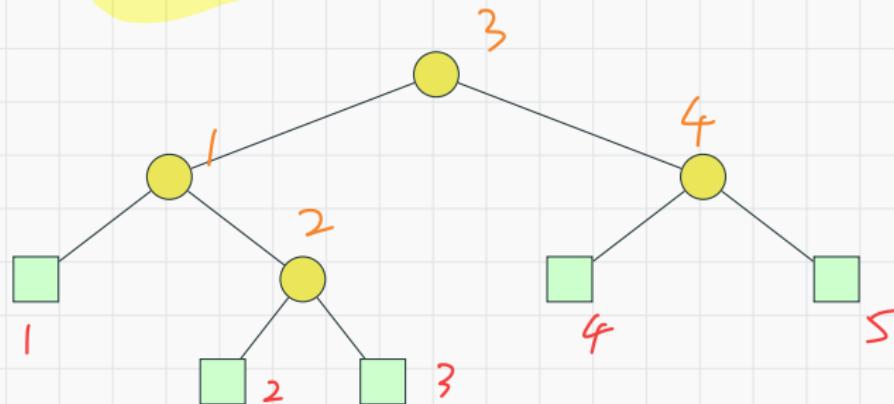


Figure 7: A full binary tree with 5 leaves and 4 internal nodes

Thus the binary tree for an optimal code for n letters contains $n - 1$ internal nodes.

HUFFMAN CODE

Huffman invented an algorithm to create an optimal code —

1. Represent each letter with a leaf node. (n leaves).

HUFFMAN CODE

Huffman invented an algorithm to create an optimal code —

1. Represent each letter with a leaf node.
2. Select two nodes with the least frequencies and connect them with a parent node whose frequency is the sum of the two.

↑
internal
node.

This step creates an internal node.

HUFFMAN CODE

Huffman invented an algorithm to create an optimal code —

1. Represent each letter with a leaf node.
2. Select two nodes with the least frequencies and connect them with a parent node whose frequency is the sum of the two.
3. replace the two selected nodes with the **parent node**.

HUFFMAN CODE

Huffman invented an algorithm to create an optimal code —

1. Represent each letter with a leaf node.
2. Select two nodes with the least frequencies and connect them with a parent node whose frequency is the sum of the two.
3. replace the two selected nodes with the parent node.
4. Repeat this until one node remains.

 The idea is that nodes with lower frequencies are merged earlier so that they are further away from the root.

 How many times do we need to run step 2 if there are n letters?

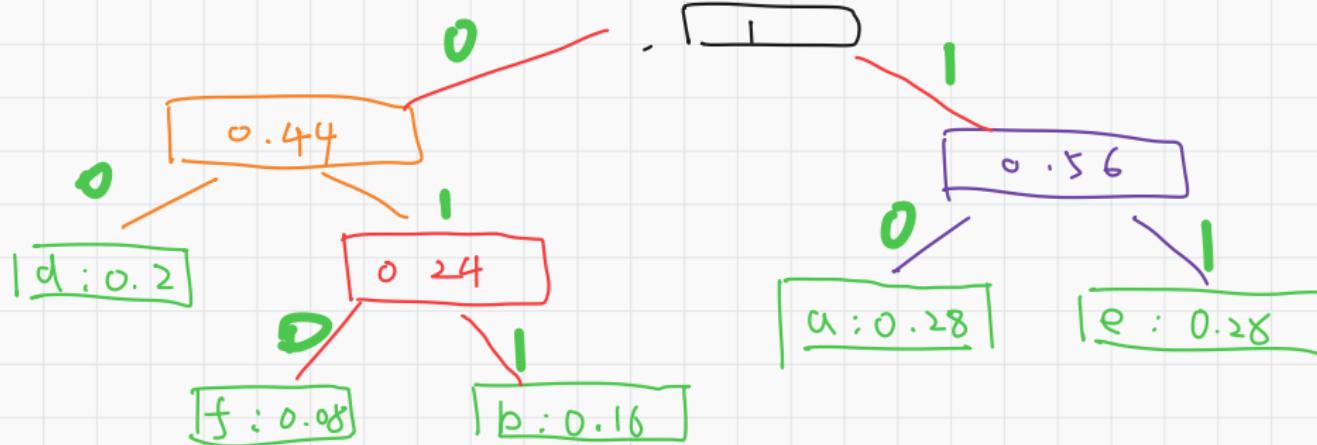
We need $n-1$ internal nodes. So we need merging $n-1$ times.

EXAMPLE OF HUFFMAN CODE

The following code is the Huffman code for the given frequency table.

Character	a	e	d	b	f
Frequency	0.28	0.28	0.2	0.16	0.08
Counts	7	7	5	4	2
Variable-Len. Codewords	10	11	00	011	010

... How can construct the Huffman code from the frequency table?



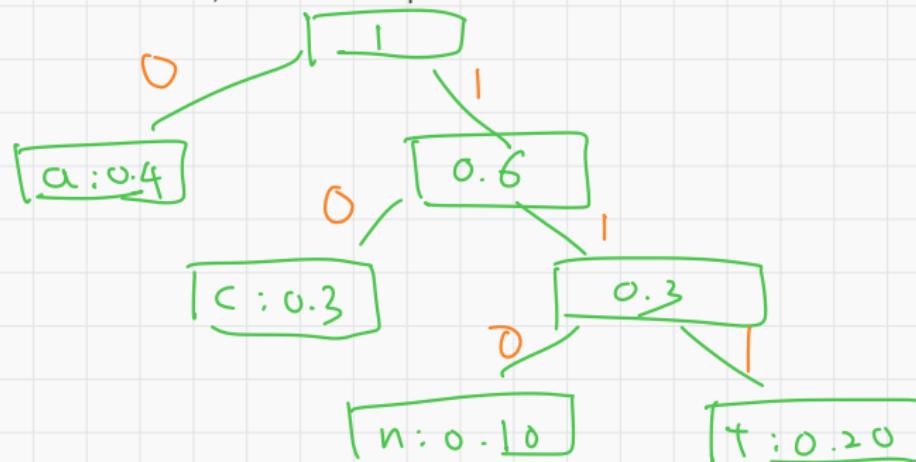
The frequency table for “a cat can act” is the following –

Character	a	c	t	n
Frequency	0.40	0.30	0.20	0.10

Code : 0 10 11 110

🎂 What is the cost of the optimal code?

💡 First find the Huffman code, then compute its cost.



MIN-PRIORITY QUEUE

To implement Huffman's algorithm *efficiently*, we need a data structure called **min-priority queue**.

Let Q denote such a queue. Then

- EXTRACT-MIN(Q) returns an object in the queue with the minimal $freq$.
- INSERT(Q, z) stores the object of z in the queue.

We assume the ⏳ complexity of each operation is $O(\log n)$.

💡 We will treat the min-priority queue as a black box.

PSEUDOCODE OF HUFFMAN'S ALGORITHM

1: **Huffman(C)**

2: $n \leftarrow |C|$ ← The number of letters.

3: $Q \leftarrow C$ ← Put all chars in min Queue.

4: **for** $i \leftarrow 1$ to $n - 1$ **do**

$\left. \begin{matrix} n-1 \\ \text{times} \end{matrix} \right\}$

5: $z \leftarrow \text{new Node}$ ← New interval node.

6: $z.\text{left} \leftarrow x \leftarrow \text{EXTRACT-MIN}(Q)$ } two nodes with min freq.

7: $z.\text{right} \leftarrow y \leftarrow \text{EXTRACT-MIN}(Q)$

8: $z.\text{freq} \leftarrow x.\text{freq} + y.\text{freq}$ $O(\log n)$

9: $\text{INSERT}(Q, z)$ ↑ Merge.
 ↑ New node

10: **return** $\text{EXTRACT-MIN}(Q)$

What is the  complexity?

$O(n \cdot \log n)$

as the root.

▷ Return the root of the tree

into the queue.

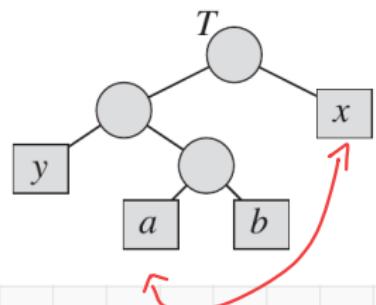
ITA 15.3 HUFFMAN CODES

CORRECTNESS OF HUFFMAN'S ALGORITHM

LEMMA 15.2

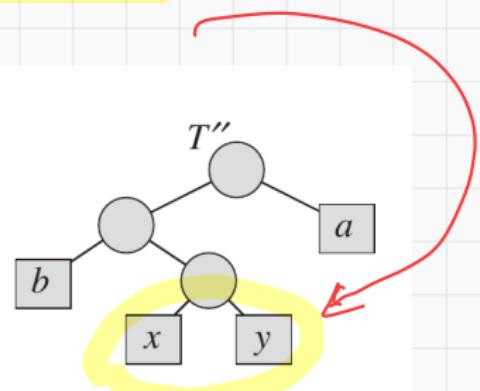
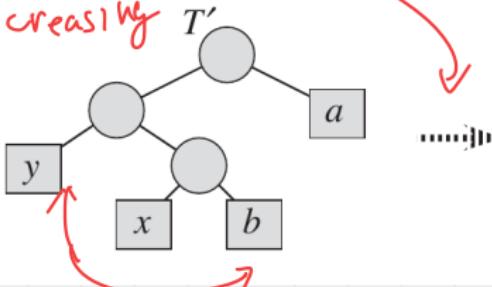
For the least frequent characters $x, y \in C$, there exists an optimal prefix code where codewords for x and y have equal length and differ only in the last bit.

Optimal.



cost not

increasing



Also Optimal.

Figure 8: A proof of Lemma 15.2

💡 This lemma means Optimal prefix-free codes have the greedy-choice property.

LEMMA 15.3

Let $x, y \in C$ be the characters with ^{least} minimum frequencies.

Let $C' = C - \{x, y\} \cup \{z\}$ where $z.freq = x.freq + y.freq$.



Replace all x and y by z .

LEMMA 15.3

Let $x, y \in C$ be the characters with minimum frequencies.

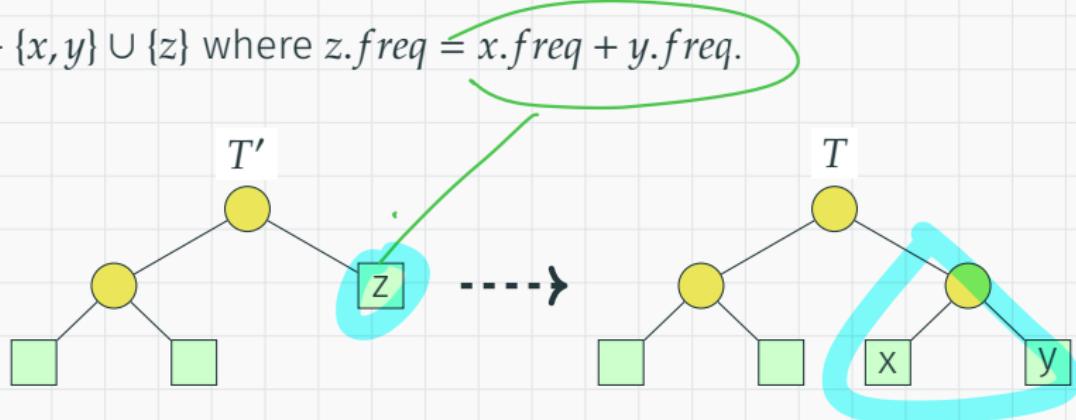
Let $C' = C - \{x, y\} \cup \{z\}$ where $z.freq = x.freq + y.freq$.



LEMMA 15.3

Let $x, y \in C$ be the characters with minimum frequencies.

Let $C' = C - \{x, y\} \cup \{z\}$ where $z.freq = x.freq + y.freq$.



For a tree T' representing an optimal prefix code for C' , the tree T formed by replacing T' 's leaf node for z with an internal node having x and y as children, yields an optimal prefix code for C .

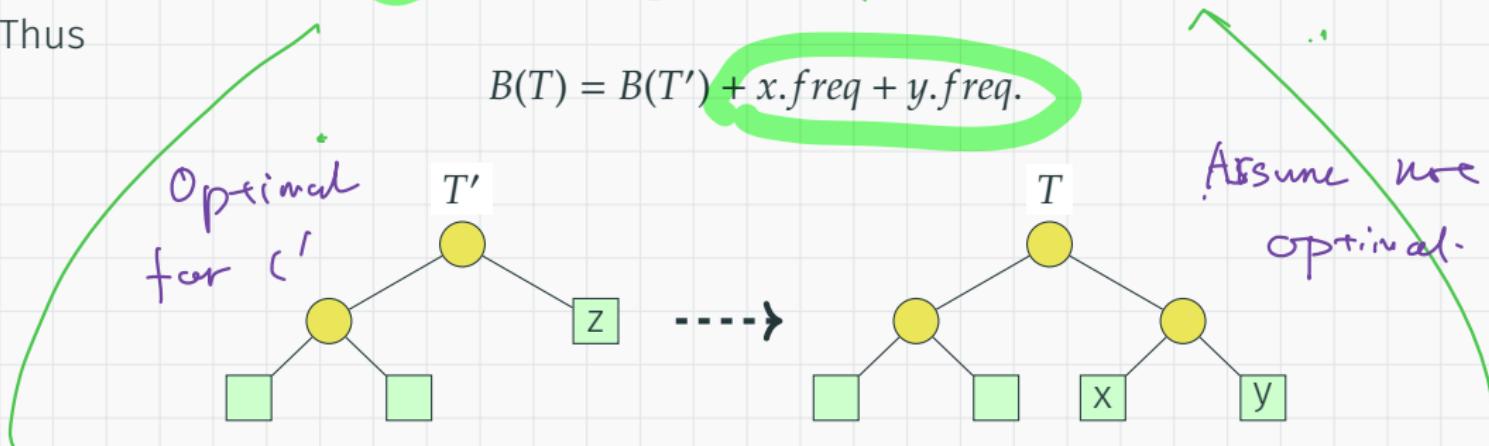
PROOF OF LEMMA 15.3

For each $c \in C - \{x, y\}$, $d_T(c) = d_{T'}(c)$. Since $d_T(x) = d_T(y) = d_{T'}(z) + 1$,

$$x.freq \cdot d_T(x) + y.freq \cdot d_T(y) = z.freq \cdot d_{T'}(z) + x.freq + y.freq.$$

Thus

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



$$= z.freq (d_{T'}(z) + 1) + y.freq (d_{T'}(z) + 1)$$

$$= (x.freq + y.freq) d_{T'}(z) + x.freq + y.freq.$$

$$= z.freq d_{T'}(z) + x.freq + y.freq.$$

PROOF OF LEMMA 15.3

Thus

$$B(T) = B(T') + x.freq + y.freq \Rightarrow \begin{aligned} & B(T') \\ &= B(T) - x.freq \\ &\quad - y.freq. \end{aligned}$$

Assume T is not optimal for C . Then there is an optimal T'' with $B(T'') < B(T)$. By Lemma 15.2, take x and y as siblings in T'' and form T''' by replacing their parent with a leaf z where $z.freq = x.freq + y.freq$. Then

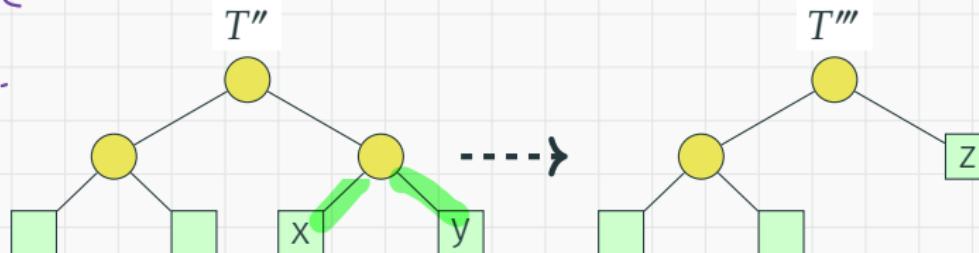
$$B(T''') = B(T'') - x.freq - y.freq < B(T'),$$

a contradiction. Therefore T is optimal for C .

$$\leftarrow B(T) - x.freq - y.freq.$$

$$= B(T).$$

Optimal
for C .



THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Base Case ($n = 2$).

When there are only two letters, the only two codes is either

- to assign one letter "o" and the other "l",
- or to swap the assignment.

Both codes thus has the same cost.

So Huffman algorithm produces an optimal prefix code.



□

THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Inductive Hypothesis — Assume that Huffman's algorithm produces an optimal prefix code for any $n - 1$ letters.



THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Inductive Hypothesis — Assume that Huffman's algorithm produces an optimal prefix code for any $n - 1$ letters.

Inductive Step —

- **Identify the Two Least Frequent Characters**

- Let x and y be the characters with the **smallest frequencies** in C .

THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Inductive Hypothesis — Assume that Huffman's algorithm produces an optimal prefix code for any $n - 1$ letters.

Inductive Step —

- **Identify the Two Least Frequent Characters**
 - Let x and y be the characters with the smallest frequencies in C .
- **Construct a Reduced Set C'**
 - Form $C' = C - \{x, y\} \cup \{z\}$, where $z.\text{freq} = x.\text{freq} + y.\text{freq}$.

THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Inductive Hypothesis — Assume that Huffman's algorithm produces an optimal prefix code for any $n - 1$ letters.

Inductive Step —

- **Identify the Two Least Frequent Characters**
 - Let x and y be the characters with the smallest frequencies in C .
- **Construct a Reduced Set C'**
 - Form $C' = C - \{x, y\} \cup \{z\}$, where $z.\text{freq} = x.\text{freq} + y.\text{freq}$.
- **Apply Inductive Hypothesis to C'**
 - By induction, Huffman's algorithm produces an optimal prefix code T' for C' .

THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Inductive Hypothesis — Assume that Huffman's algorithm produces an optimal prefix code for any $n - 1$ letters.

Inductive Step —

- **Identify the Two Least Frequent Characters**
 - Let x and y be the characters with the smallest frequencies in C .
- **Construct a Reduced Set C'**
 - Form $C' = C - \{x, y\} \cup \{z\}$, where $z.\text{freq} = x.\text{freq} + y.\text{freq}$.
- **Apply Inductive Hypothesis to C'**
 - By induction, Huffman's algorithm produces an optimal prefix code T' for C' .
- **Expand T' to Obtain T for C**
 - Replace the leaf node for z in T' with an internal node having x and y as children (per Lemma 15.3).

THEOREM 15.4

Huffman's algorithm produces an optimal prefix code.

Proof.

Inductive Hypothesis — Assume that Huffman's algorithm produces an optimal prefix code for any $n - 1$ letters.

Inductive Step —

- **Identify the Two Least Frequent Characters**
 - Let x and y be the characters with the smallest frequencies in C .
- **Construct a Reduced Set C'**
 - Form $C' = C - \{x, y\} \cup \{z\}$, where $z.\text{freq} = x.\text{freq} + y.\text{freq}$.
- **Apply Inductive Hypothesis to C'**
 - By induction, Huffman's algorithm produces an optimal prefix code T' for C' .
- **Expand T' to Obtain T for C**
 - Replace the leaf node for z in T' with an internal node having x and y as children (per Lemma 15.3).
- **Conclude Optimality** $\textcolor{green}{T}$ is also exactly what Huffman's algorithm produces.

HUFFMAN CODES WITH FIBONACCI FREQUENCIES

Symbol	Huffman Code	Frequency
a	1	1
b	0	1



Do you see the pattern?

HUFFMAN CODES WITH FIBONACCI FREQUENCIES

Symbol	Huffman Code	Frequency
a	11	1
b	10	1
c	0	2



Do you see the pattern?

HUFFMAN CODES WITH FIBONACCI FREQUENCIES

Symbol	Huffman Code	Frequency
a	111	1
b	110	1
c	10	2
d	0	3



Do you see the pattern?

HUFFMAN CODES WITH FIBONACCI FREQUENCIES

Symbol	Huffman Code	Frequency
a	1111	1
b	1110	1
c	110	2
d	10	3
e	0	5



Do you see the pattern?

HUFFMAN CODES WITH FIBONACCI FREQUENCIES

Symbol	Huffman Code	Frequency
a	11111	1
b	11110	1
c	1110	2
d	110	3
e	10	5
f	0	8



Do you see the pattern?

HUFFMAN CODES WITH FIBONACCI FREQUENCIES

Symbol	Huffman Code	Frequency
a	111111	1
b	111110	1
c	11110	2
d	1110	3
e	110	5
f	10	8
g	0	13



Do you see the pattern?

😜 BEYOND HUFFMAN CODING: NOTABLE NEW DIRECTIONS IN COMPRESSION (2025-2026)



Figure 9: Silicon Valley is an American comedy television series about company trying to find new compression algorithms

Learned compression. Neural codecs learn transforms and models. Entropy coding still finishes the job.

Learned compression. Neural codecs learn transforms and models. Entropy coding still finishes the job.

Generative compression. The decoder guesses fine details. Smaller files, less exact images.

Learned compression. Neural codecs learn transforms and models. Entropy coding still finishes the job.

Generative compression. The decoder guesses fine details. Smaller files, less exact images.

Compression for machines (VCM). Bitstreams keep features for AI tasks. Human visuals may look worse.

Learned compression. Neural codecs learn transforms and models. Entropy coding still finishes the job.

Generative compression. The decoder guesses fine details. Smaller files, less exact images.

Compression for machines (VCM). Bitstreams keep features for AI tasks. Human visuals may look worse.

AV2 video standard. AV2 follows AV1. It targets smaller files.

Learned compression. Neural codecs learn transforms and models. Entropy coding still finishes the job.

Generative compression. The decoder guesses fine details. Smaller files, less exact images.

Compression for machines (VCM). Bitstreams keep features for AI tasks. Human visuals may look worse.

AV2 video standard. AV2 follows AV1. It targets smaller files.

JPEG XL renaissance. JPEG XL compresses old JPEGs. It may return in browsers.

WHAT ARE YOUR MAIN TAKEAWAYS TODAY? ANY QUESTIONS?

