

UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

COMP 550 Algorithm and Analysis

Greedy Algorithms (Huffman Encoding)

Based on CLRS Sec. 15.3

Some slides are adapted from ones by prior instructors Prof. Plaisted and Prof. Osborne

Data Compression

Lossy

Big file \rightarrow Small file \rightarrow Imperfect copy



- Example: JPEG, MP3
- Idea: Targeted data deletion.
 - Delete parts of data that people won't notice.

Lossless

Big file \rightarrow Small file \rightarrow Perfect copy



- Example: ZIP, PNG, JBIG.
- How to do this?

Fixed-Length Encoding

- Let's consider compression of text files
- How is text commonly represented?
 - Fixed length encoding (e.g., ASCII, Unicode)
 - 8 bits in ASCII
- Assuming ASCII, how many bits represent this text (ignore whitespace)?

she sells seashells by the seashore

- Number of characters: 30
- Bits required: $30 \times 8 = 240$

Fixed-Length Encoding

• Assuming ASCII, how many bits represent this text (ignore whitespace)?

she sells seashells by the seashore

- Number of characters: 30
- Bits required: $30 \times 8 = 240$
- There are 10 different characters in this text
 - At least 4 bits per character are required under fixed-length encoding scheme
 - Bits required: $30 \times 4 = 120$

COMP550@UNC

Variable-Length Encoding

- Drawbacks of fixed-length codes
 - Wasted space
 - Same number of bits used to represent all characters
- Potential solution: use variable-length codes
 - Variable number of bits when frequency of occurrence is known
 - Short codes for characters that occur frequently

	Letter	%				
	E	13				
	Т	9.1				
	A	8.2				
	0	7.5				
	I	7				
•	N	6.7				
	all others					

English letter frequencies (Wikipedia)

Variable-Length Encoding

she sells seashells by the seashore



Letter	%
E	13
Т	9.1
A	8.2
0	7.5
I	7
N	6.7
all others	

Variable-Length Encoding

she sells seashells by the seashore



Number of bits = 86

Variable-Length Encoding

- One issue: Where a character ends and another begins?
 - Not a problem under fixed-length encoding



Prefix-Free Codes

• No codeword is a prefix of another codeword

			Letter	Freq	Bit code	Bit count	Bits required
			S	8	10	2	8×2=16
		1	E		00	2	7×2=14
	No other codeword		Н	4	011	3	4×3=12
	starts with 10		L	4	110	3	4x3=12
Code: 0110011011011111 HELLO		А	2	0101	4	2×4=8	
		R	1	0100	4	1×4=4	
		В	1	11100	5	1×5=5	
		У	1	11101	5	1×5=5	
		Т	1	11110	5	1×5=5	
		0	1	11111	5	1~5-5	

Prefix-Free Codes

• The following codes do not satisfy prefix-free property

Letter	Bit code
S	0
E	1
Н	01
L	10

- Decode the code 1010
 - Can be anyone among "ESES, LL, LES, ESL"

Prefix-Free Codes



Optimal Encoding

- <u>Problem</u>: Determine optimal prefix-free codes for a given text
 - Optimal = Minimize the number of bits to store the text
 - Minimize: \sum frequecy of a char \times Number of bits for the char
- <u>Key property:</u>
 - Less frequent characters are further down the tree

- 1. Tabulate character frequencies
- 2. Each char: frequency tuple forms a single-node tree



















3. Merge the two trees with the smallest roots. The new root is the sum of

the two frequencies below it. Repeat until there's one tree.



3. Merge the two trees with the smallest roots. The new root is the sum of

the two frequencies below it. Repeat until there's one tree.



4. Label left edges "0" and right edges "1"



Codes are given by the path from root to leaf



- <u>Input</u>: A set C of n characters each with an attribute freq
- <u>Output</u>: A tree corresponding to optimal code



Correctness

- To show that Huffman encoding builds an optimal prefix-free code:
 - The greedy choice of merging two trees with the smallest root satisfies the greedy-choice property
 - Constructing optimal prefix-free codes has the optimal substructure property

Greedy-Choice Property

There is an optimal solution containing the greedy choice (lowest frequency characters are siblings at the maximum depth)

Proof (sketch) by exchange argument:

- Take a tree T representing an arbitrary optimal prefix-free code
- Transform the tree T to a tree T' that has the greedy choice made by Huffman algorithm
- Show that T' is no worse than T



Optimal Substructure Property

Making the least frequency characters siblings and solving the new subproblem results in an optimal prefix-free code.

<u>Goal</u>: Replace the least frequency character x and y and a new node z with

z.freq = x.freq + y.freq

Optimal solution of problem = optimal solution of subproblem + (x and y as children of z) <u>Proof sketch</u>: Assume that

Optimal solution of problem \neq optimal solution of subproblem + (x and y as children of z)

- By the greedy-choice property, we know there is an optimal solution *OPT* where x and y are siblings.
- Removing x and y from OPT and replacing their parents in OPT by z results in a better solution for the subproblem, contradiction

Encoding & Decoding

- Encoding:
 - Create a hashmap or dictionary or array (indexed by characters)
 to store bit code of each character
 - Scan each character in the text, look for its code, and concatenate
- Decoding:
 - Go left/right the tree representing prefix-free code according to the current bit
 - If leaf is reached, concatenate that character to current text

Thank You!