

## Data Compression: Huffman Coding

10.1 in Weiss (p.389)

1

## Why compress files?

2

## Why compress files?

- For long term storage (disc space is limited)
- For transferring files over the internet (bigger files take longer)
- A smaller file more likely to fit in memory/cache

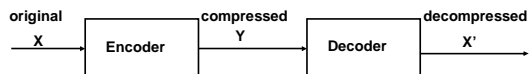
3

## What is a file?

- C++ program code
- Executable program
- Email - text
- HTML document
- Pictures (lossy); JPEG
- Video (lossy); MPEG
- Audio (lossy); MP3

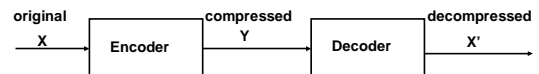
4

## Data Compression



5

## Data Compression



- **Lossless** compression  $X = X'$
- **Lossy** compression  $X \neq X'$
- **Compression Ratio**  $|X|/|Y|$ 
  - Where  $|X|$  is the # of bits in X.

6

## Lossy Compression

- Some data is lost, but not too much.

### Standards:

- JPEG (Joint Photographic Experts Group) – stills
- MPEG (Motion Picture Experts Group)
  - Audio and video
- MP3 (MPEG-1, Layer 3)

7

## Lossless Compression

- No data is lost.

### Standards:

- Gzip, Unix compress, zip, GIF, Morse code
- Examples:
  - Run-length Encoding (RLE)
  - Huffman Coding

8

## RLE

- Idea: Compactly represent long ‘runs’ of the same character
- “aaarrrrr!” as ‘a’x3 ‘r’x5 then ‘!’

9

## RLE

- Idea: Compactly represent long ‘runs’ of the same character
- “aaarrrrr!” as ‘a’x3 ‘r’x5 then ‘!’
- Say...
  - Replace all ‘runs’ of the same character by 2 characters: the 1) character and 2) the length
  - ‘bee’ becomes ‘b’,1,’e’,2

10

## RLE

- Idea: Compactly represent long ‘runs’ of the same character
- “aaarrrrr!” as ‘a’x3 ‘r’x5 then ‘!’
- Say...
  - Replace all ‘runs’ of the same character by 2 characters: the 1) character and 2) the length
  - ‘bee’ becomes ‘b’,1,’e’,2
  - When is this good?
  - When is this really bad?

11

## Another idea: Use fewer bits per character

**ASCII** = fixed 8 bits per character

**Example:** “hello there”

- 11 characters \* 8 bits = 88 bits

Can we encode this message using fewer bits?

12

## Another idea: Use fewer bits per character

ASCII = fixed 8 bits per character

**Example:** "hello there"

– 11 characters \* 8 bits = 88 bits

Can we encode this message using fewer bits?

- We could look JUST at the message
- there are only 6 possible characters + one space = 7 things; only need 3 bits
- Encode: aabdcaa = could do as 16 bits (each character = 2 bits each)
- Huffman can do as 14 bits

13

## Huffman Coding

- Uses *frequencies* of symbols in a string to build a **prefix code**.
- **Prefix Code** – no code in our encoding is a prefix of another code.

Letter	code
a	0
b	100
c	101
d	11

14

## Huffman Coding

- Uses *frequencies* of symbols in a string to build a **prefix code**.
- **Prefix Code** – no code in our encoding is a prefix of another code.

Letter	code
a	0
b	100
c	101
d	11

15

## Huffman Coding

- Uses *frequencies* of symbols in a string to build a **prefix code**.
- **Prefix Code** – no code in our encoding is a prefix of another code.

Letter	code
a	0
b	100
c	101
d	11

16

## Decoding a Prefix Code

Loop

start at root of tree

loop

if bit read = 1 then go right

else, go left

until node is a leaf

Report character found!

Until end of the message

17

## Decode: 11100010100110

Letter	code
a	0
b	100
c	101
d	11

18

## Decode: 11100010100110

Letter	code
a	0
b	100
c	101
d	11

19

## Huffman Trees

Cost of a Huffman Tree containing n symbols

$$C(T) = p_1 * r_1 + p_2 * r_2 + p_3 * r_3 + \dots + p_n * r_n$$

Where:

$p_i$  = the probability that a symbol occurs

$r_i$  = the length of the path from the root to the node

20

## Example Cost

Letter	Frequency	code
a	.50	0
b	.125	100
c	.125	101
d	.25	11

Cost: 1.75

21

## Constructing a tree

- Determine frequency of each letter/symbol
- Place each as an unconnected leaf node
- Repeatedly merge two nodes with lowest frequency into one node with sum of frequencies
- Huffman Coding is optimal\*

22

## Constructing a tree example

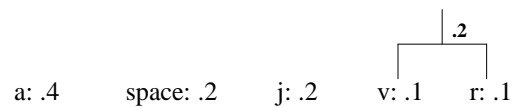
- Encode "a java jar"
- 4 a's, 2 spaces, 2 j's, 1 v, 1 r; 10 total

a: .4    space: .2    j: .2    v: .1    r: .1

23

## Constructing a tree example

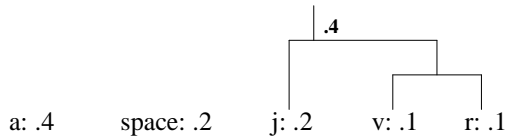
- Encode "a java jar"
- 4 a's, 2 spaces, 2 j's, 1 v, 1 r; 10 total



24

### Constructing a tree example

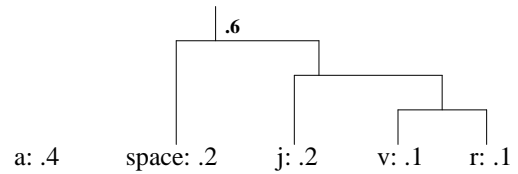
- Encode “a java jar”
- 4 a’s, 2 spaces, 2 j’s, 1 v, 1 r; 10 total



25

### Constructing a tree example

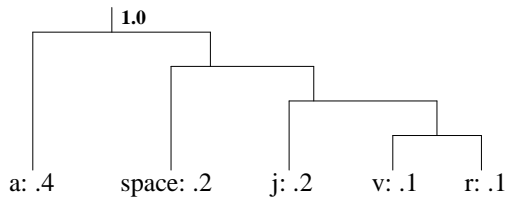
- Encode “a java jar”
- 4 a’s, 2 spaces, 2 j’s, 1 v, 1 r; 10 total



26

### Constructing a tree example

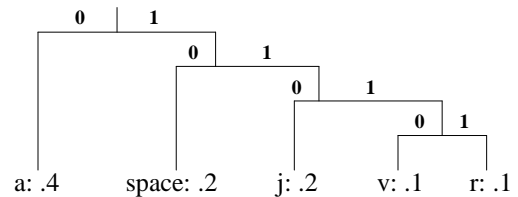
- Encode “a java jar”
- 4 a’s, 2 spaces, 2 j’s, 1 v, 1 r; 10 total



27

### Constructing a tree example

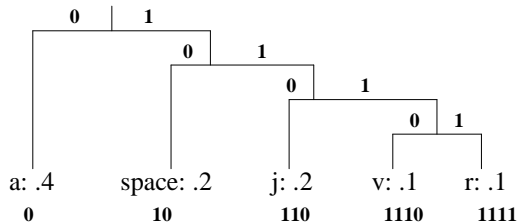
- Encode “a java jar”
- 4 a’s, 2 spaces, 2 j’s, 1 v, 1 r; 10 total



28

### Constructing a tree example

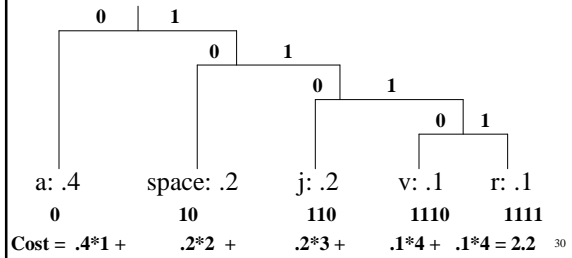
- Encode “a java jar”
- 4 a’s, 2 spaces, 2 j’s, 1 v, 1 r; 10 total



29

### Constructing a tree example

- Encode “a java jar”
- 4 a’s, 2 spaces, 2 j’s, 1 v, 1 r; 10 total



30

## Run-time?

- To decode an encoded message length  $n$ :

31

## Run-time?

- To decode an encoded message length  $n$ :  $O(n)$
- To encode message length  $n$ , with  $c$  possible characters

32

## Run-time?

- To decode an encoded message length  $n$ :  $O(n)$
- To encode message length  $n$ , with  $c$  possible characters
  - Count frequencies:
  - Build tree:
  - Encode:

33

## Run-time?

- To decode an encoded message length  $n$ :  $O(n)$
- To encode message length  $n$ , with  $c$  possible characters
  - Count frequencies:  $O(n)$
  - Build tree:  $O(\log c)$  (with priority queue)
  - Encode:  $O(n)$

34