

Greedy algorithms 2

David Kauchak
cs302
Spring 2012



Knapsack problems: Greedy or not?

- **0-1 Knapsack** – A thief robbing a store finds n items worth v_1, v_2, \dots, v_n dollars and weight w_1, w_2, \dots, w_n pounds, where v_i and w_i are integers. The thief can carry at most W pounds in the knapsack. Which items should the thief take if he wants to maximize value.
- **Fractional knapsack problem** – Same as above, but the thief happens to be at the bulk section of the store and can carry fractional portions of the items. For example, the thief could take 20% of item i for a weight of $0.2w_i$ and a value of $0.2v_i$.



Data compression

- Given a file containing some data of a fixed alphabet Σ (e.g. A, B, C, D), we would like to pick a binary character code that minimizes the number of bits required to represent the data.

ACADAADB ...

➔

0010100100100 ...

minimize the size of the encoded file



Compression algorithms

General purpose [10]

- Huffman coding (HFC) – a simple scheme that provides good compression of data consisting only of one of the seven values
- LZ77 (LZ) – LZ77, LZ78, LZSS, LZJ, LZH, LZMA, LZMA2, LZMA3 – used by ZIP, WinRAR, and various other file formats
- LZ78 (LZ) – used by ZIP, WinRAR, and various other file formats
- LZSS (LZ) – LZSS (Lempel-Ziv-Storer-Szymanski) – used by ZIP, WinRAR, and various other file formats
- LZMA (LZ) – LZMA (Lempel-Ziv-Markov chain algorithm) – used by 7-Zip, xz, and other programs; higher compression than LZSS/LZ78
- LZMA2 (LZ) – LZMA2 (Lempel-Ziv-Markov chain algorithm) – designed for compression/decompression speed at the expense of compression ratio
- LZMA3 (LZ) – LZMA3 (Lempel-Ziv-Markov chain algorithm) – a combination of statistical method and dictionary-based method; faster compression ratio than LZMA2

Audio [10]

- Free Lossless Audio Codec – FLAC
- Apple Lossless – ALAC (Apple Lossless Audio Codec)
- MP3 – Lossy
- Adaptive Transform Acoustic Coding – ATAC
- Audio Lossless Coding – also known as MERGA-ALC
- MPEG-4 AAC – also known as HEAAC
- Short Range Transform – SRT
- Dolby TrueHD
- VORBIS Music Audio
- Meridian Lossless Packing – MLP
- Monkey Audio – Monkey's Audio APE
- DASH/MP3
- Original Sound Quality – OSQ
- RealAudio – RealAudio Lossless
- TrueAudio – TTA
- True Audio Lossless
- WavPack – WavPack Lossless
- Xiph.Org – Xiph.Org Music Lossless

Graphics [10]

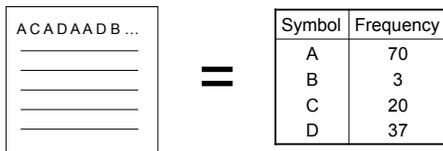
- LZW – Lossless LZW compression of ASCII/UTF images
- JPEG – Standard lossy compression of RGB images
- JPEG-LS – Standard lossless compression of images
- JPEG-2000 – Standard lossless compression method, authored by Scott Koonin, Prof. Sun Diego State University
- JBIG – Progressive Scalable File Interchange for lossy compression
- PNG – Portable Network Graphics
- PPM – Pseudo Planar File Format
- GIF – Graphics Interchange Format
- GIF87 – GIF87a (GIF89) – GIF87a (GIF89) – GIF87a (GIF89) – GIF87a (GIF89)
- GIF89 – GIF89 (GIF89) – GIF89 (GIF89) – GIF89 (GIF89)

http://en.wikipedia.org/wiki/Lossless_data_compression



Simplifying assumption: frequency only

Assume that we only have character frequency information for a file



Fixed length code

Use $\lceil \log_2 |\Sigma| \rceil$ bits for each character

- A =
- B =
- C =
- D =

Fixed length code

Use $\lceil \log_2 |\Sigma| \rceil$ bits for each character

- A = 00 $2 \times 70 +$
 - B = 01 $2 \times 3 +$
 - C = 10 $2 \times 20 +$
 - D = 11 $2 \times 37 =$
- 260 bits

Symbol	Frequency
A	70
B	3
C	20
D	37

How many bits to encode the file?

Fixed length code

Use $\lceil \log_2 |\Sigma| \rceil$ bits for each character

- A = 00 $2 \times 70 +$
 - B = 01 $2 \times 3 +$
 - C = 10 $2 \times 20 +$
 - D = 11 $2 \times 37 =$
- 260 bits

Symbol	Frequency
A	70
B	3
C	20
D	37

Can we do better?

Variable length code

What about:

A = 0 $1 \times 70 +$
 B = 01 $2 \times 3 +$
 C = 10 $2 \times 20 +$
 D = 1 $1 \times 37 =$
 153 bits

Symbol	Frequency
A	70
B	3
C	20
D	37

How many bits to encode the file?

Decoding a file

A = 0 010100011010
 B = 01
 C = 10
 D = 1

What characters does this sequence represent?

Decoding a file

A = 0 010100011010
 B = 01 $\underbrace{\hspace{1.5em}}$
 C = 10 **A D or B?**
 D = 1

What characters does this sequence represent?

Variable length code

What about:

A = 0 $1 \times 70 +$
 B = 100 $3 \times 3 +$
 C = 101 $3 \times 20 +$
 D = 11 $2 \times 37 =$

Symbol	Frequency
A	70
B	3
C	20
D	37

213 bits
 (18% reduction)

How many bits to encode the file?

Prefix codes

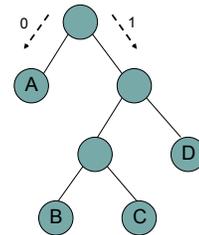
A prefix code is a set of codes where no codeword is a **prefix** of some other codeword

A = 0	A = 0
B = 01	B = 100
C = 10	C = 101
D = 1	D = 11

Prefix tree

We can encode a prefix code using a **full** binary tree where each child represents an encoding of a symbol

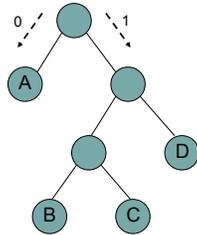
A = 0
B = 100
C = 101
D = 11



Decoding using a prefix tree

- To decode, we traverse the graph until a leaf node is reached and output the symbol

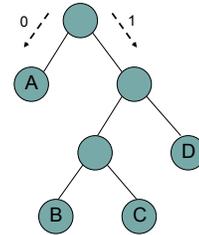
A = 0
B = 100
C = 101
D = 11



Decoding using a prefix tree

- Traverse the graph until a leaf node is reached and output the symbol

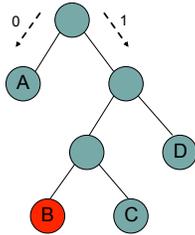
1000111010100



Decoding using a prefix tree

- Traverse the graph until a leaf node is reached and output the symbol

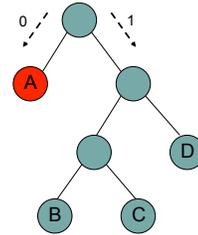
1000111010100
B



Decoding using a prefix tree

- Traverse the graph until a leaf node is reached and output the symbol

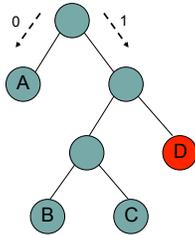
1000111010100
B A



Decoding using a prefix tree

- Traverse the graph until a leaf node is reached and output the symbol

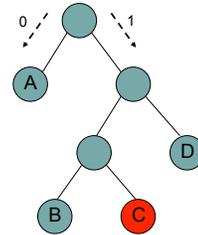
1000111010100
B A D



Decoding using a prefix tree

- Traverse the graph until a leaf node is reached and output the symbol

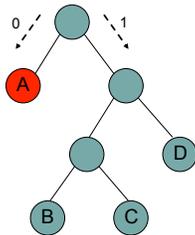
1000111010100
B A D C



Decoding using a prefix tree

- Traverse the graph until a leaf node is reached and output the symbol

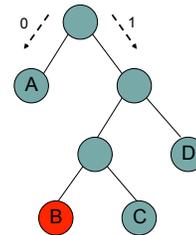
1000111010100
 B A D C A



Decoding using a prefix tree

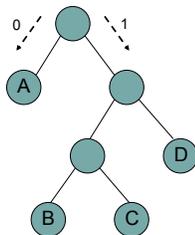
- Traverse the graph until a leaf node is reached and output the symbol

1000111010100
 B A D C A B



Determining the cost of a file

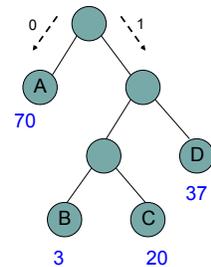
Symbol	Frequency
A	70
B	3
C	20
D	37



Determining the cost of a file

Symbol	Frequency
A	70
B	3
C	20
D	37

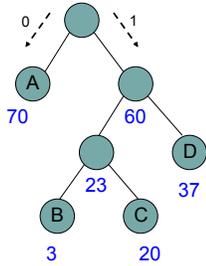
$$\text{cost}(T) = \sum_{i=1}^n f_i \text{depth}(i)$$



Determining the cost of a file

Symbol	Frequency
A	70
B	3
C	20
D	37

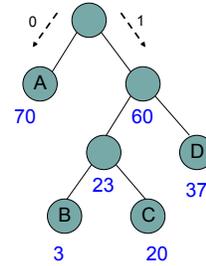
What if we label the internal nodes with the sum of the children?



Determining the cost of a file

Symbol	Frequency
A	70
B	3
C	20
D	37

Cost is equal to the sum of the internal nodes and the leaf nodes



Determining the cost of a file

As we move down the tree, one bit gets read for every nonroot node

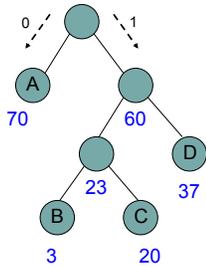
70 times we see a 0 by itself

60 times we see a prefix that starts with a 1

of those, 37 times we see an additional 1

the remaining 23 times we see an additional 0

of these, 20 times we see a last 1 and 3 times a last 0



A greedy algorithm?

Given file frequencies, can we come up with a prefix-free encoding (i.e. build a prefix tree) that minimizes the number of bits?

```

HUFFMAN(F)
1  Q ← MAKEHEAP(F)
2  for i ← 1 to |Q| - 1
3      allocate a new node z
4      left[z] ← x ← EXTRACTMIN(Q)
5      right[z] ← y ← EXTRACTMIN(Q)
6      f[z] ← f[x] + f[y]
7      INSERT(Q, z)
8  return EXTRACTMIN(Q)
    
```

```

HUFFMAN(F)
1 Q ← MAKEHEAP(F)
2 for i ← 1 to |Q| - 1
3   allocate a new node z
4   left[z] ← x ← EXTRACTMIN(Q)
5   right[z] ← y ← EXTRACTMIN(Q)
6   f[z] ← f[x] + f[y]
7   INSERT(Q,z)
8 return EXTRACTMIN(Q)
    
```

Symbol	Frequency
A	70
B	3
C	20
D	37

Heap

```

HUFFMAN(F)
1 Q ← MAKEHEAP(F)
2 for i ← 1 to |Q| - 1
3   allocate a new node z
4   left[z] ← x ← EXTRACTMIN(Q)
5   right[z] ← y ← EXTRACTMIN(Q)
6   f[z] ← f[x] + f[y]
7   INSERT(Q,z)
8 return EXTRACTMIN(Q)
    
```

Symbol	Frequency
A	70
B	3
C	20
D	37

Heap

B 3
C 20
D 37
A 70

```

HUFFMAN(F)
1 Q ← MAKEHEAP(F)
2 for i ← 1 to |Q| - 1
3   allocate a new node z
4   left[z] ← x ← EXTRACTMIN(Q)
5   right[z] ← y ← EXTRACTMIN(Q)
6   f[z] ← f[x] + f[y]
7   INSERT(Q,z)
8 return EXTRACTMIN(Q)
    
```

Symbol	Frequency
A	70
B	3
C	20
D	37

Heap

merging with this node will incur an additional cost of 23

BC 23
D 37
A 70

```

HUFFMAN(F)
1 Q ← MAKEHEAP(F)
2 for i ← 1 to |Q| - 1
3   allocate a new node z
4   left[z] ← x ← EXTRACTMIN(Q)
5   right[z] ← y ← EXTRACTMIN(Q)
6   f[z] ← f[x] + f[y]
7   INSERT(Q,z)
8 return EXTRACTMIN(Q)
    
```

Symbol	Frequency
A	70
B	3
C	20
D	37

Heap

BCD 60
A 70

```

HUFFMAN(F)
1  Q ← MAKEHEAP(F)
2  for i ← 1 to |Q| - 1
3      allocate a new node z
4      left[z] ← x ← EXTRACTMIN(Q)
5      right[z] ← y ← EXTRACTMIN(Q)
6      f[z] ← f[x] + f[y]
7      INSERT(Q, z)
8  return EXTRACTMIN(Q)
    
```

Symbol	Frequency
A	70
B	3
C	20
D	37

Heap

ABCD 130

Is it correct?

- The algorithm selects the symbols with the two smallest frequencies first (call them f_1 and f_2)

Is it correct?

- The algorithm selects the symbols with the two smallest frequencies first (call them f_1 and f_2)
- Consider a tree that did not do this:

Is it correct?

- The algorithm selects the symbols with the two smallest frequencies first (call them f_1 and f_2)
- Consider a tree that did not do this:

$$\text{cost}(T) = \sum_{i=1}^n f_i \text{depth}(i)$$

- frequencies don't change
- cost will **decrease** since $f_1 < f_i$

contradiction

Runtime?

```

HUFFMAN( $F$ )
1  $Q \leftarrow \text{MAKEHEAP}(F)$ 
2 for  $i \leftarrow 1$  to  $|Q| - 1$ 
3   allocate a new node  $z$ 
4    $left[z] \leftarrow x \leftarrow \text{EXTRACTMIN}(Q)$ 
5    $right[z] \leftarrow y \leftarrow \text{EXTRACTMIN}(Q)$ 
6    $f[z] \leftarrow f[x] + f[y]$ 
7    $\text{INSERT}(Q, z)$ 
8 return  $\text{EXTRACTMIN}(Q)$ 

```

1 call to MakeHeap

2(n-1) calls ExtractMin

n-1 calls Insert

$O(n \log n)$

Non-optimal greedy algorithms

- All the greedy algorithms we've looked at today give the optimal answer
- Some of the most common greedy algorithms generate good, but non-optimal solutions
 - set cover
 - clustering
 - hill-climbing
 - relaxation