


# Greedy algorithms


---

David Kauchak  
cs302  
Spring 2013



## Administrative

?




## Greedy algorithms

What is a greedy algorithm?

Algorithm that makes a local decision with the goal of creating a globally optimal solution

Method for solving problems where optimal solutions can be defined in terms of optimal solutions to sub-problems


What does this mean? Where have we seen this before?



## Greedy vs. divide and conquer


### Divide and conquer

To solve the general problem:




↓

Break into sum number of sub problems, solve:



then possibly do a little work



## Greedy vs. divide and conquer



Divide and conquer

To solve the general problem:



The solution to the general problem is solved with respect to solutions to sub-problems!

## Greedy vs. divide and conquer



Greedy

To solve the general problem:



Pick a locally optimal solution and repeat



## Greedy vs. divide and conquer



Greedy

To solve the general problem:



The solution to the general problem is solved with respect to solutions to sub-problems!

Slightly different than divide and conquer

## Horn formula



A horn formula is a set of implications and negative clauses:

$$\Rightarrow x \quad x \wedge u \Rightarrow z$$

$$\Rightarrow y \quad \bar{x} \vee \bar{y} \vee \bar{z}$$

**Goal**

Given a horn formula, determine if the formula is satisfiable, i.e. an assignment of true/false to the variables that is consistent with all of the implications/causes

$$\Rightarrow x \quad x \wedge u \Rightarrow z$$

$$\Rightarrow y \quad \bar{x} \vee \bar{y} \vee \bar{z}$$

u	x	y	z
0	1	1	0

**A greedy solution?**

$$\Rightarrow x \quad x \wedge z \Rightarrow w \quad w \wedge y \wedge z \Rightarrow x$$

$$x \Rightarrow y \quad x \wedge y \Rightarrow w \quad \bar{w} \vee \bar{x} \vee \bar{y}$$

$$w \ 0$$

$$x \ 0$$

$$y \ 0$$

$$z \ 0$$

**A greedy solution?**

$$\Rightarrow x \quad x \wedge z \Rightarrow w \quad w \wedge y \wedge z \Rightarrow x$$

$$x \Rightarrow y \quad x \wedge y \Rightarrow w \quad \bar{w} \vee \bar{x} \vee \bar{y}$$

$$w \ 0$$

$$x \ 1$$

$$y \ 0$$

$$z \ 0$$

**A greedy solution?**

$$\Rightarrow x \quad x \wedge z \Rightarrow w \quad w \wedge y \wedge z \Rightarrow x$$

$$x \Rightarrow y \quad x \wedge y \Rightarrow w \quad \bar{w} \vee \bar{x} \vee \bar{y}$$

$$w \ 0$$

$$x \ 1$$

$$y \ 1$$

$$z \ 0$$

### A greedy solution?

$$\begin{array}{lll} \Rightarrow x & x \wedge z \Rightarrow w & w \wedge y \wedge z \Rightarrow x \\ x \Rightarrow y & \boxed{x \wedge y \Rightarrow w} & \bar{w} \vee \bar{x} \vee \bar{y} \end{array}$$

w 1  
x 1  
y 1  
z 0

### A greedy solution?

$$\begin{array}{lll} \Rightarrow x & x \wedge z \Rightarrow w & w \wedge y \wedge z \Rightarrow x \\ x \Rightarrow y & x \wedge y \Rightarrow w & \boxed{\bar{w} \vee \bar{x} \vee \bar{y}} \end{array}$$

w 1  
x 1  
y 1  
z 0

not satisfiable

### A greedy solution

```

HORN(H)
1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true

```

### A greedy solution

```

HORN(H)
1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true

```

set all variables of the implications of the form " $\Rightarrow x$ " to true

## A greedy solution

HORN( $H$ )

```

1 set all variables to false
2 for all implications  $i$ 
3   if EMPTY(LHS( $i$ ))
4     RHS( $i$ )  $\leftarrow$  true
5 changed  $\leftarrow$  true
6 while changed
7   changed  $\leftarrow$  false
8   for all implications  $i$ 
9     if LHS( $i$ ) = true and !RHS( $i$ ) = true
10      RHS( $i$ )  $\leftarrow$  true
11      changed = true
12 for all negative clauses  $c$ 
13   if  $c$  = false
14     return false
15 return true

```

if the all variables of the lhs of an implication are true, then set the rhs variable to true

## A greedy solution

HORN( $H$ )

```

1 set all variables to false
2 for all implications  $i$ 
3   if EMPTY(LHS( $i$ ))
4     RHS( $i$ )  $\leftarrow$  true
5 changed  $\leftarrow$  true
6 while changed
7   changed  $\leftarrow$  false
8   for all implications  $i$ 
9     if LHS( $i$ ) = true and !RHS( $i$ ) = true
10      RHS( $i$ )  $\leftarrow$  true
11      changed = true
12 for all negative clauses  $c$ 
13   if  $c$  = false
14     return false
15 return true

```

see if all of the negative clauses are satisfied

## Correctness of greedy solution

Two parts:

- If our algorithm returns an assignment, is it a valid assignment?
- If our algorithm does not return an assignment, does an assignment exist?

## Correctness of greedy solution

If our algorithm returns an assignment, is it a valid assignment?

HORN( $H$ )

```

1 set all variables to false
2 for all implications  $i$ 
3   if EMPTY(LHS( $i$ ))
4     RHS( $i$ )  $\leftarrow$  true
5 changed  $\leftarrow$  true
6 while changed
7   changed  $\leftarrow$  false
8   for all implications  $i$ 
9     if LHS( $i$ ) = true and !RHS( $i$ ) = true
10      RHS( $i$ )  $\leftarrow$  true
11      changed = true
12 for all negative clauses  $c$ 
13   if  $c$  = false
14     return false
15 return true

```

## Correctness of greedy solution

If our algorithm returns an assignment, is it a valid assignment?

```

HORN(H)
1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true
  
```

explicitly check all negative clauses

## Correctness of greedy solution

If our algorithm returns an assignment, is it a valid assignment?

```

HORN(H)
1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true
  
```

don't stop until all implications with all lhs elements true have rhs true

## Correctness of greedy solution

If our algorithm does not return an assignment, does an assignment exist?

```

HORN(H)
1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true
  
```

Our algorithm is "stingy". It only sets those variables that **have to be true**. All others remain false.

## Running time?

```

HORN(H)
1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true
  
```

?  
n = number of variables  
m = number of formulas

## Running time?

HORN(*H*)

```

1 set all variables to false
2 for all implications i
3   if EMPTY(LHS(i))
4     RHS(i) ← true
5 changed ← true
6 while changed
7   changed ← false
8   for all implications i
9     if LHS(i) = true and !RHS(i) = true
10      RHS(i) ← true
11      changed = true
12 for all negative clauses c
13   if c = false
14     return false
15 return true
    
```

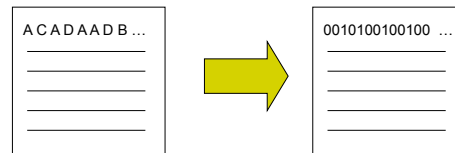
$O(nm)$

*n* = number of variables  
*m* = number of formulas

## Data compression

Given a file containing some data of a fixed alphabet  $\Sigma$  (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.

minimize the size of the encoded file

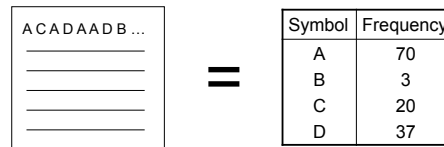


## Compression algorithms

- General purpose**
- Run-length encoding (RLE) - a simple scheme that provides good compression of data consisting lots of runs of the same value
  - Lempel-Ziv 1978 (LZ78), Lempel-Ziv-Stolzi (LZ78) - used by GIF images and compress among many other applications
  - GZIP (LZ77) - widely used, ZIP archive format (LZ77) and as part of the compressed process of Portable Network Graphics (PNG), Portable Pixmap Format (PPM), TIFF, ISO
  - LZSS - using the Lempel-Ziv scheme, this provides almost the highest compression than LZ78/LZ77
  - Lempel-Ziv - Welch (LZW) - designed for compression/compression based on the expense of compression ratio
  - Statistical Lempel-Ziv - a combination of statistical method and dictionary based method, better compression with than using single method
- Audio**
- Free Lossless Audio Codec - FLAC
  - Apple Lossless - ALAC (Apple Lossless Audio Codec)
  - MP3 - standard
  - Adaptive Transform Acoustic Coding - ATAC
  - Fast Lossless Coding - also known as MPEG-4 AAC
  - MPEG-4 AAC - also known as AAC
  - DASH (Dynamic Adaptive Streaming over HTTP)
  - Dolby TrueHD
  - DTS-HD Master Audio
  - Microsoft Lossless Packing - MLP
  - Microsoft Audio - Microsoft Audio (MMA)
  - OggVorbis
  - Original Sound Quality - OSQ
  - RealAudio - RealAudio Lossless
  - TrueHD - True Audio Lossless
  - TrueHD - RealAudio Lossless
  - TrueHD - RealAudio Lossless
  - TrueHD - RealAudio Lossless
- Graphics**
- LZW - Lossless RLE compression of images (GIF images)
  - JPEG - standard in lossy compression of images
  - JPEG-LS - Lossless lossless compression standard
  - JPEG-LS - Lossless lossless compression method, as proven by Sunil Kumar, Prof. Sanjiv Kumar, Prof. Sanjiv Kumar (State University)
  - JPEG-LS - Lossless lossless compression method, as proven by Sunil Kumar, Prof. Sanjiv Kumar, Prof. Sanjiv Kumar (State University)
  - PNG - Portable Network Graphics
  - TIFF - Tagged Image File Format
  - GIF - Graphics Interchange Format
  - GIF - Graphics Interchange Format
  - GIF - Graphics Interchange Format
- [http://en.wikipedia.org/wiki/Lossless\\_data\\_compression](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

Can we do better?

Symbol	Frequency
A	70
B	3
C	20
D	37



## 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            {  
 C = 10            A D or B?  
 D = 1

What characters does this sequence represent?

## Variable length code

What about:

A = 0  
 B = 100    Is it decodeable?  
 C = 101  
 D = 11

Symbol	Frequency
A	70
B	3
C	20
D	37

## Variable length code

What about:

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

213 bits  
 (18% reduction)

Symbol	Frequency
A	70
B	3
C	20
D	37

How many bits to encode the file?

### Prefix codes

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

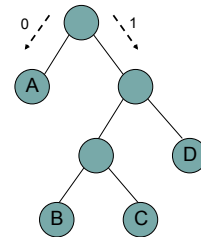
A = 0	A = 0
B = 01	B = 100
C = 10	C = 101
D = 1	D = 11



### Prefix tree

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

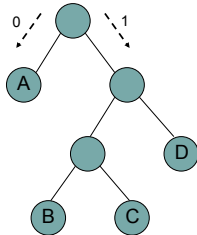
A = 0
B = 100
C = 101
D = 11



### Decoding using a prefix tree

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

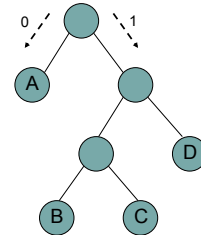
A = 0
B = 100
C = 101
D = 11



### Decoding using a prefix tree

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

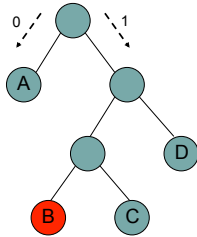
1000111010100



### Decoding using a prefix tree

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

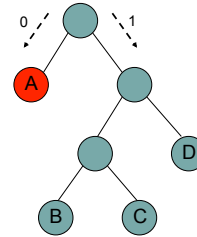
1000111010100  
B



### Decoding using a prefix tree

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

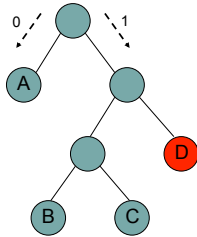
1000111010100  
B A



### Decoding using a prefix tree

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

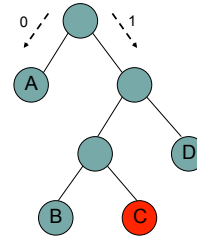
1000111010100  
B A D



### Decoding using a prefix tree

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

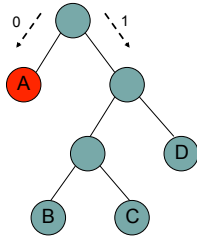
1000111010100  
B A D C



### Decoding using a prefix tree

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

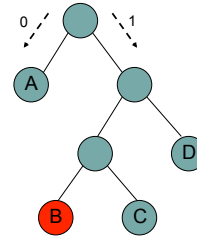
1000111010100  
B A D C A



### Decoding using a prefix tree

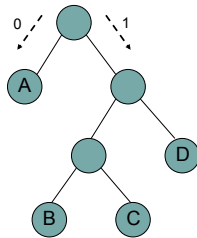
Traverse the graph until a leaf node is reached and output the symbol

1000111010100  
B A D C A B



### Determining the cost of a file

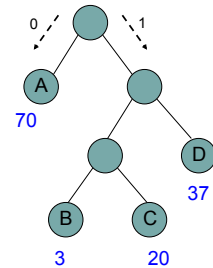
Symbol	Frequency
A	70
B	3
C	20
D	37



### Determining the cost of a file

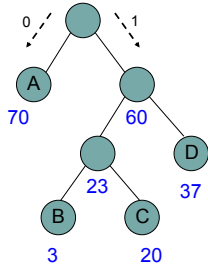
Symbol	Frequency
A	70
B	3
C	20
D	37

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



### Determining the cost of a file

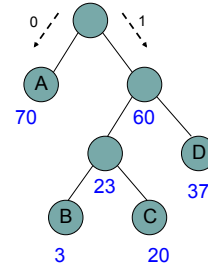
Symbol	Frequency
A	70
B	3
C	20
D	37



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

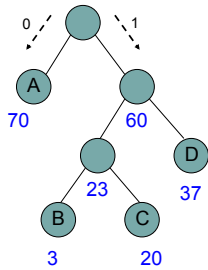
70 times we see a 0 by itself

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

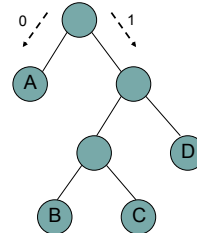
of these, 20 times we see a last 1 and 3 times a last 0



### 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?

Symbol	Frequency
A	70
B	3
C	20
D	37



## 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

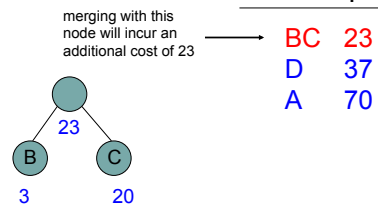
Heap

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



```

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  $f_2$ )

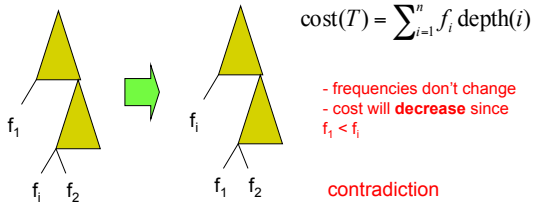
Consider a tree that did not do this (proof by contradiction):

Is it optimal?

## 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:



## Runtime?

```

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)

```

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 so far give the optimal answer

Some of the most common greedy algorithms generate good, but non-optimal solutions

- set cover
- clustering
- hill-climbing
- relaxation

## 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  $0.2v_i$ .