

CS3000
6/2 - Mon.

Admin

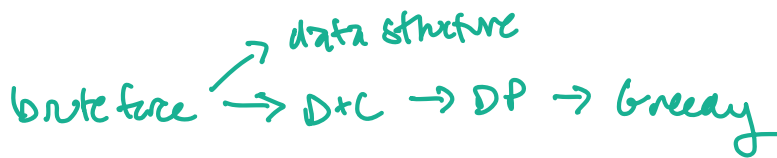
- HW4 due Thurs 9pm
- APPE due 11:30 tmcw
- Thu 6/12 exam #21,
 - ↳ practice problems at 6/5,
 - review in rec on 6/10

Agenda

1. File compression
2. Greedy Huffman
3. APPE

0. Greedy Strategy

can we do better?



What is greedy?

- makes optimal choice at each step
- make it easy to make the optimal choice
- need to prove correctness

What kinds of problems?

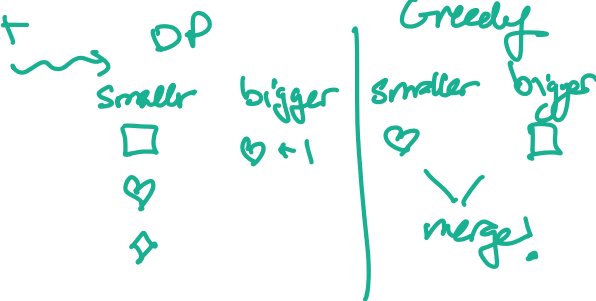
- optimization problems
- optimal substructure

Pros/Cons vs. DP?

- easier to implement
- careful about correctness
- (usually) faster

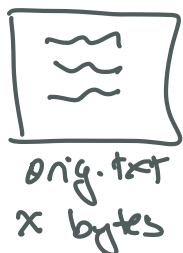
Both DP, Greedy:

- need optimal substructure
- make local decisions at each step



1. File Compression

- text file of characters
 - ↳ flat file
 - ↳ one byte per character

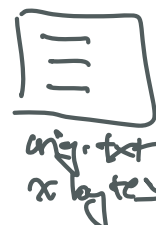


— compress —
mapping
char: encoded



compressed.txt
y bytes
• $y < x$
• lossless

— decompress —
tree
encoded: char



ex) File: Larry IS Awesome 1000 times

total space: $(5 + 2 + 7) \times 1000 = 14,000 = 14 \text{ KB}$

compress: Larry - L

IS - I

Awesome - A

LIA 1000 times

$(1 + 1 + 1) \times 1000 = 3,000 = 3 \text{ KB}$

valid! But, what if another word starts with L?
What about optimization?

valid solution: compress
decompress
lossless

optimal solution: minimize space!

Huffman: compress characters, not words

• char = 1 byte (8 bits)

• if we can replace with ≤ 7 bits that's a win!

ex) File: AABBBBBB C

orig. txt

space $9 \times 8 = 72$ bits

three chars to encode: A, B, C

• one in one bit ~ most frequent (B)

• two in two bits ~ A, C

Building the tree:

- highest freq character = closest to root

- characters are leaves

- 0 on left edges

- 1 on right edges

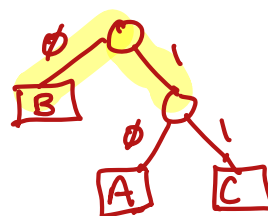
mapping
char \rightarrow encoded:

B: 0

A: 10

C: 11

Huffman Tree



decompress [ex]

100010

A B B A

Huffman gives a prefix-free code

• no way to encode a char that is a prefix for something else

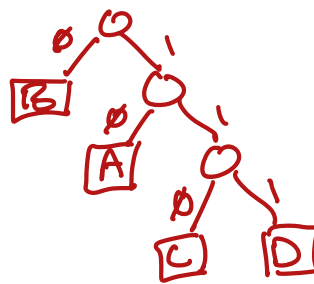
\rightarrow 01 bad !!

• only one way to decode any bit sequence

Where would D go?

- assume D.freq = 1

010111 → BAD!



Huffman Algorithm

- goal: build the tree!
- optimization goal: tree creates the smallest encoding

For the algo:

- C, set of chars
- |C| # of characters
- priority queue Q
 - extract min
 - insert
- every char c has c.freq

Chars to use:

2	b	c	d	e	f
45	13	12	16	9	5

Qs to consider:

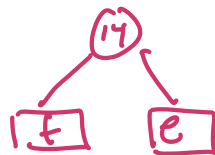
- keys in nodes vs. leaves?
- optimal solution?
 - ↳ How many bits do we start with? 900
- what is greedy choice? → 2 lowest frequencies
- what is in Q each step?

HUFFMAN(C)

```

1  n = |C|
2  Q = C
3  for i = 1 to n - 1
4      allocate a new node z
5      x = EXTRACT-MIN(Q)
6      y = EXTRACT-MIN(Q)
7      z.left = x
8      z.right = y
9      z.freq = x.freq + y.freq
10     INSERT(Q, z)
11  return EXTRACT-MIN(Q)
    
```

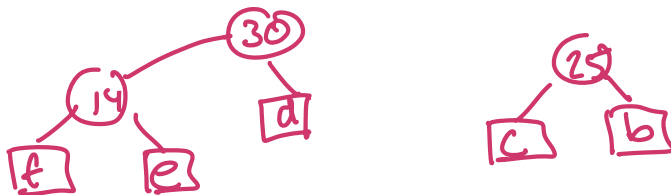
↳ Q: 45 13 12 16 9 5



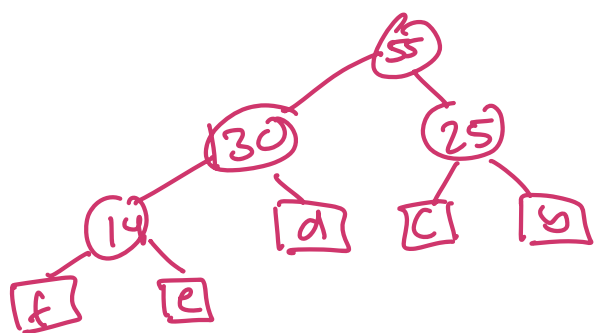
↳ Q: 45 13 12 16 14



↳ Q: 45 25 16 14

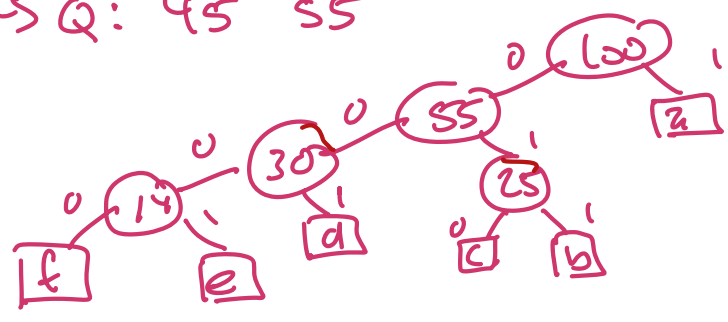


↳ Q: 45 25 30



0010110011
D B D A

↳ Q: 45 55

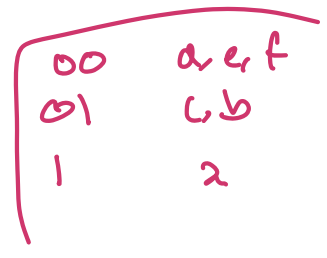


bits now!
 a = 1 (1) d = 3 (001)
 b = 3 (011) c = 4 (0001)
 c = 3 (010) e = 4 (0000)

compressed bits: $1 \cdot 45 + 3 \cdot 13 + 3 \cdot 12 + 3 \cdot 16 + 4 \cdot 9 + 4$

= 224 bits !!

not prefix free:
 0110 101
 0101
 0100



APPG