

MORE GREEDY ALGORITHMS

Huffman Codes - optimal prefix code

ASCII is a fixed length code

$$\begin{array}{ll} A = 41_{16} = 0100 0001 & a = 61_{16} = 0110 0001 \\ B = 42_{16} = 0100 0010 & b = 62_{16} = 0110 0010 \\ \vdots & \vdots \\ Z = 5A_{16} = 0101 1010 & z = 7A_{16} = 0111 1010 \end{array}$$

Each character uses 8 bits, no matter how often the character is used

Variable-length codes

shorter codes for frequently used characters
longer codes for infrequent characters

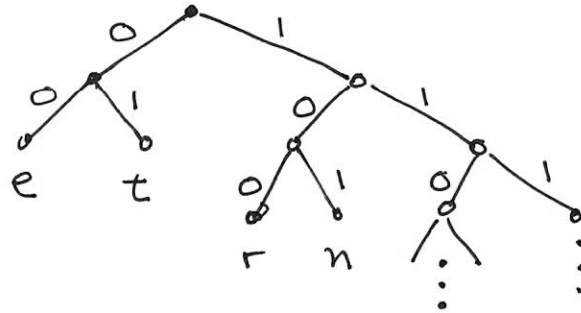
Example: VCR-plus ↪ who uses VCRs anymore?

Question: how to decode? 1000 0101 01
↪ where do letters start?

Prefix-free code = prefix code ↪ boneless chicken = boxed chicken

No character has an encoding that
is the prefix of another character's encoding

Represent prefix codes as a binary tree:



e=00
t=01
r=100
n=101

1000 0101 01 01
r e n t

Problem: Given an "alphabet" of characters and a frequency for each character, what is the optimal prefix code?

Minimize

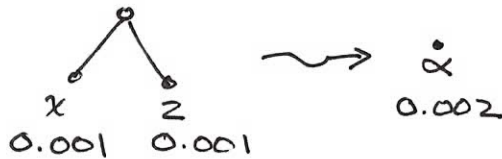
$$B(T) = \sum_{x \in C} f(x) d_T(x)$$

Annotations:

- frequency of x (points to $f(x)$)
- depth of x in tree = length of encoding of x (points to $d_T(x)$)
- char (points to x in the summation)
- char set (points to C in the summation)
- tree for prefix code (points to the summation symbol)
- expected # of bits (points to $B(T)$)

Greedy algorithm: Huffman codes

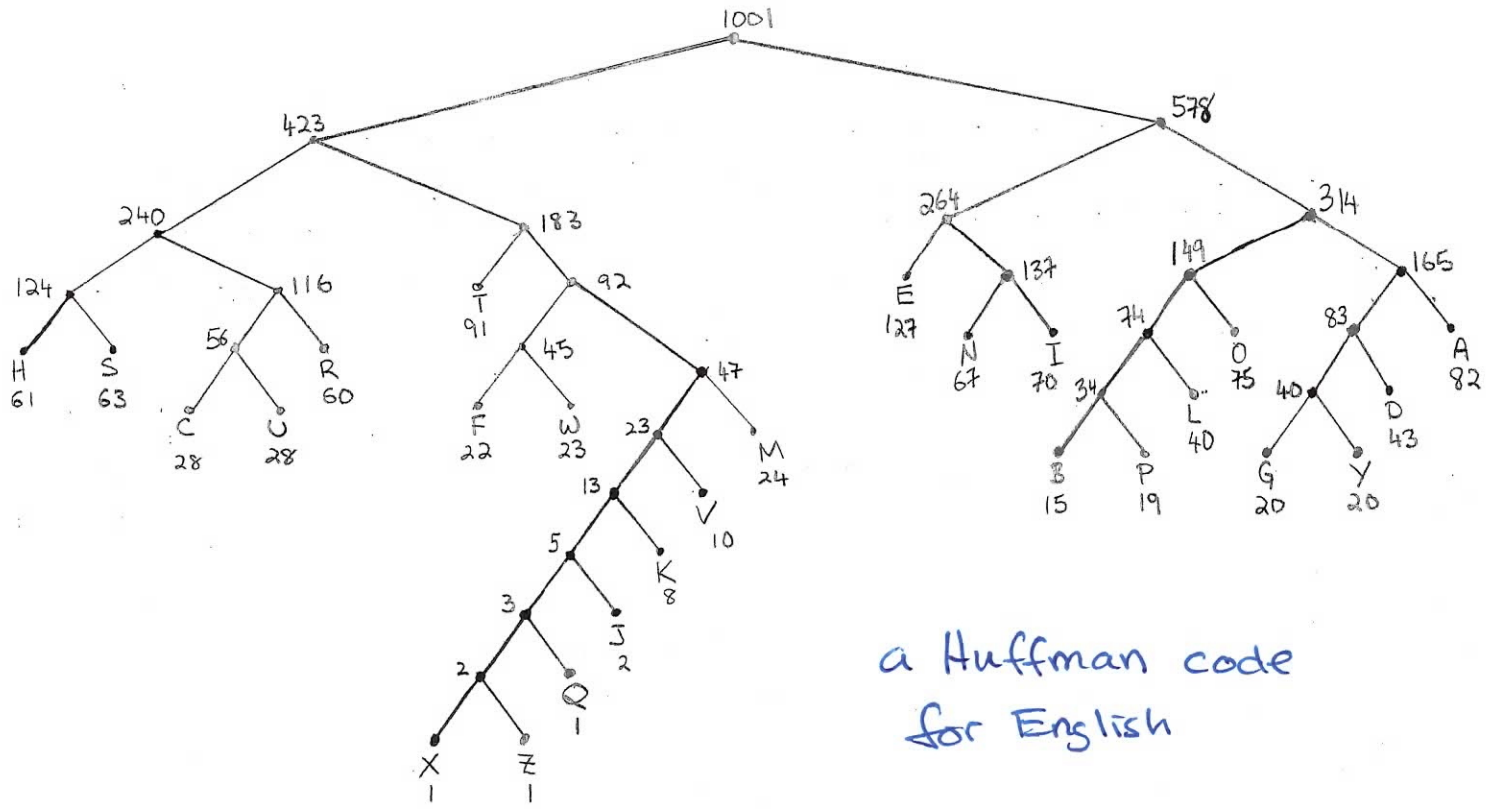
- ① Sort characters by frequency
- ② Take 2 characters with least frequency: x & z .
(These will be placed at the bottom.)
- ③ Remove x & z from the alphabet and replace them with a "fake" character with



$$C' = (C - \{x, z\}) \cup \{\alpha\}$$

letter	probability	letter	probability
A	.082	N	.067
B	.015	O	.075
C	.028	P	.019
D	.043	Q	.001
E	.127	R	.060
F	.022	S	.063
G	.020	T	.091
H	.061	U	.028
I	.070	V	.010
J	.002	W	.023
K	.008	X	.001
L	.040	Y	.020
M	.024	Z	.001

letter frequency table for English



a Huffman code
for English

Running Time:

Use a heap to store characters by frequency.
Smaller frequencies on top.

Each iteration = 2 ExtractMin + 1 Insert

Build Heap		$\Theta(n)$
Extract Min x2	} n times	$2 \times n \times \Theta(\log n)$
Insert		$n \times \Theta(\log n)$
		<hr/>
		$\Theta(n \log n)$

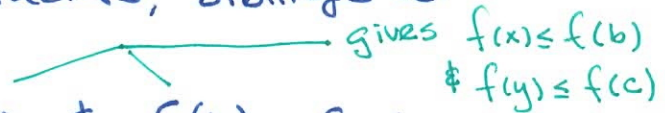
Huffman Codes are optimum, but why?

Claim #1: Suppose T is a binary tree for some prefix code.

Let $x \neq y$ be two characters with lowest frequency.

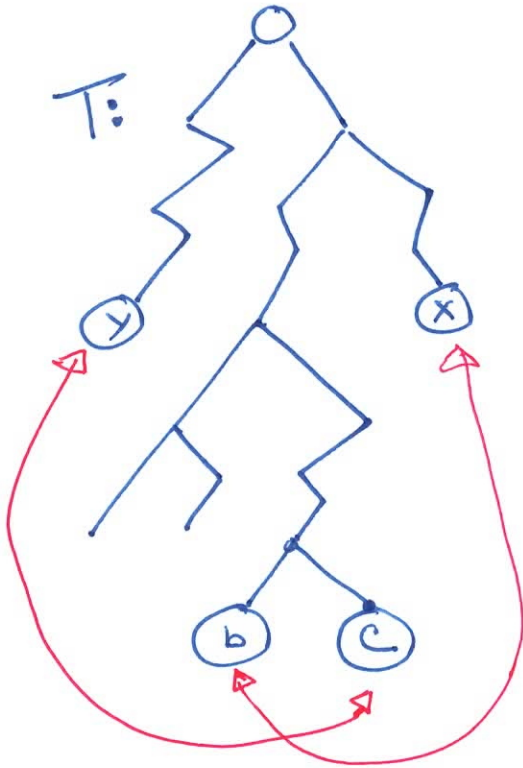
Let $b \neq c$ be two characters, siblings at maximum depth.

Suppose that $f(x) \leq f(y) \neq f(b) \leq f(c)$.

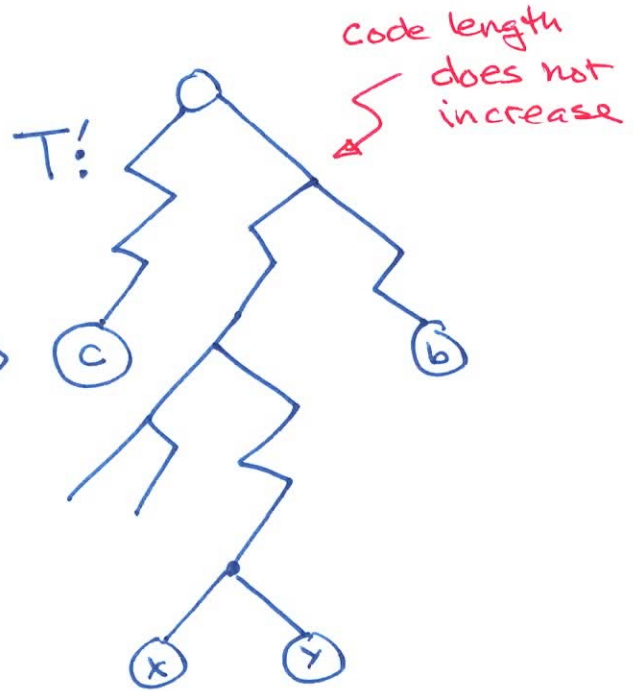


Then, swapping x and b and swapping y and c will result in T' with no greater expected code length.

Claim #1:



Swap
⇒



Proof of Claim #1: (by calculation)

$$\begin{aligned} B(T') &= B(T) - \underbrace{f(x) \cdot d_T(x)}_{\substack{\text{removed } x \\ \text{at depth} \\ d_T(x)}} + \underbrace{f(x) \cdot d_{T'}(x)}_{\substack{\text{added } x \\ \text{at new depth}}} \leftarrow = d_T(b) \\ &\quad - f(y) \cdot d_T(y) + f(y) \cdot d_{T'}(y) \leftarrow = d_T(c) \\ &\quad - f(b) \cdot d_T(b) + f(b) d_{T'}(b) \\ &\quad - f(c) d_T(c) + f(c) d_{T'}(c) \end{aligned}$$

$$\begin{aligned} &= B(T) + f(x) \cdot [d_T(b) - d_T(x)] \\ &\quad + f(y) \cdot [d_T(c) - d_T(y)] \\ &\quad + f(b) \cdot [d_T(x) - d_T(b)] \\ &\quad + f(c) \cdot [d_T(y) - d_T(c)] \end{aligned}$$

lowest freq. (green arrow pointing to $f(b)$)

negations (red arrows pointing to the brackets)

$$\begin{aligned} &= B(T) + \underbrace{[f(b) - f(x)]}_{\geq 0} \underbrace{[d_T(x) - d_T(b)]}_{\leq 0} + \underbrace{[f(c) - f(y)]}_{\geq 0} \underbrace{[d_T(y) - d_T(c)]}_{\leq 0} \\ &\leq B(T) \end{aligned}$$

max depth (green arrow pointing to $d_T(b)$)

□

Claim #2: Let T be a binary tree for a prefix code.

Suppose $x \neq y$ are sibling leaves.

Let T' be obtained from T by replacing the subtree



with α

where α is assigned the frequency of $x \neq y$ combined: $f(\alpha) = f(x) + f(y)$.

Then, $B(T') = B(T) - f(x) - f(y)$

Proof of Claim #2: (by calculation)

$$\begin{aligned} B(T') &= B(T) - f(x) d_T(x) - f(y) d_T(y) + f(\alpha) d_{T'}(\alpha) \\ &= B(T) - f(x) d_T(x) - f(y) d_T(x) + f(\alpha) [d_T(x) - 1] \\ &= B(T) - [f(x) + f(y)] d_T(x) + f(\alpha) d_T(x) - f(\alpha) \\ &= B(T) - \cancel{f(\alpha) \cdot d_T(x)} + \cancel{f(\alpha) d_T(x)} - f(\alpha) \\ &= B(T) - f(x) - f(y) \end{aligned}$$

↙ does not depend on
the depth of x & y .

☒

Theorem: Huffman codes are optimum.

Proof by induction: Induction on # of characters.

Induction Hypothesis: Huffman codes are optimum } call this
for alphabets with n characters. } $P(n)$.

Base Case: $n=0$. Trivial.

Induction Case: Prove $P(n) \Rightarrow P(n+1)$

Suppose not. Let X be a prefix code st.

$$B(X) < B(T)$$

where T is the tree for a Huffman code.

Let $x \neq y$ be two characters at maximum depth of T . Then, $x \neq y$ have the lowest frequency.

By Claim #1, we can swap $x \neq y$ to maximum depth in X and obtain

$$B(x') \leq B(x)$$

By Claim #2, we can replace $x \neq y$ (now guaranteed to be siblings) with α in X' and in T , and obtain X'' and T'

$$B(X'') = B(X') - f(x) - f(y)$$

$$B(T') = B(T) - f(x) - f(y)$$

Then, $B(X') \leq B(X) < B(T)$

$$\Rightarrow B(X') - f(x) - f(y) < B(T) - f(x) - f(y)$$

$$\Rightarrow B(X'') < B(T')$$

But, T' is the Huffman code for the n character alphabet $C' = (C - \{x, y\}) \cup \{a\}$.

Furthermore, X'' is also a prefix code for C' .

Then, Huffman codes are not optimum for C' which has n characters. This contradicts $P(n)$.



Tape Archive Problem:

n files stored on magnetic tape.

Order matters:



must scan through
files 1 thru $i-1$ to
reach file i

Some files accessed more frequently.

Two ideas:

- ① Put shorter files in front
- ② Put more frequently used files in front.

Problems:

- ① Lots of short files that are not used often
- ② Very long files used most frequently.

Solution: use frequency / length ratio.

i^{th} file has length l_i & frequency P_i

Assume files are already sorted by P_i/l_i ratio

$$P_1/l_1 \geq P_2/l_2 \geq P_3/l_3 \geq \dots \geq P_n/l_n$$

Show that this minimize expected access time

$$\sum_{i=1}^n P_i (\text{access time for file } i) = \sum_{i=1}^n P_i \left(\underbrace{\sum_{j=1}^i l_j}_{\text{read over all files before file } i \text{ and file } i \text{ itself.}} \right)$$

read over all files before
file i and file i itself.

Usual argument: if some ordering X does not place file 1 at the beginning, it could have.

WARNING: Cannot swap to the front.

EXAMPLE: 3 files

$$p_1 = 0.80 \quad l_1 = 80$$

$$p_1/l_1 = 0.01$$

$$p_2 = 0.10 \quad l_2 = 100$$

$$p_2/l_2 = 0.001$$

$$p_3 = 0.10 \quad l_3 = 9$$

$$p_3/l_3 = 0.0111\dots$$

order:

$$1\ 2\ 3 \quad 0.8(80) + 0.10(80+100) + 0.10(80+100+9) = 100.9$$

$$1\ 3\ 2 \quad 0.8(80) + 0.10(89) + 0.10(189) = 91.8$$

$$2\ 1\ 3 \quad 0.10(100) + 0.80(180) + 0.10(189) = 172.9$$

$$2\ 3\ 1 \quad 0.10(100) + 0.10(109) + 0.80(189) = 172.1$$

$$3\ 1\ 2 \quad 0.10(9) + 0.80(89) + 0.10(189) = 91 \quad \leftarrow \text{opt.}$$

$$3\ 2\ 1 \quad 0.10(9) + 0.10(109) + 0.80(189) = 163$$

swapping
3 to the front
increases
expected
access
time

But we can swap adjacent files

S : file file file ... file $_i$ file $_j$... file ...

and $P_j/l_j > P_i/l_i$. swap file i & file j .

S' = file . file file ... file $_j$ file $_i$... file ...

Expected
Access Time

access time
not changed

access time
not changed

$$EAT(S') - EAT(S) = P_i l_j - P_j l_i$$

must read
file j to get
to file i

no longer need
to read file i to get
to file j

But, $P_j/l_j > P_i/l_i \Rightarrow P_i l_j - P_j l_i < 0$. So, $EAT(S') < EAT(S)$.

Let G be the ordering produced by greedy:

file 1 . file 2 ... file n .


How do we make S look like G

by swapping adjacent files ???

Answer: Bubble Sort

Let $S^{(t)}$ be the ordering obtained from Bubble Sorting S according to largest P_i/l_i ratio at the t^{th} step of Bubble Sort.

Then, $EAT(S) \geq EAT(S^{(1)}) \geq EAT(S^{(2)}) \geq \dots \geq EAT(S^{(m)})$

Sorted! 

Proof that Greedy is optimum:

Greedy
Solution

Suppose some solution S has $EAT(S) < EAT(G)$.

Then, Bubble Sort S by P_i/e_i order.

$$EAT(G) > EAT(S) \geq EAT(S^{(1)}) \geq \dots \geq EAT(S^{(m)}) = EAT(G)$$

... method

same ordering

∴ $EAT(G) > EAT(G)$
a contradiction.

