# Lecture #10

# 2 More Greedy Algorithms

- Huffman Encoding
    - used to compress data
    - gzip, jpeg, mp3,...
    - Intro to Entropy

- Horn Clauses
    - special case of satisfiability of a Boolean expression
      $(x$ or $\bar{y})$ and $(\bar{z}$ or $w)$ and ...
    - used in logic programming (Prolog,...)

1

# Data Compression Problem - 1

Given string of chars:  A B A C C D B B...
from alphabet $\{A, B, C, D\}$ how many bits
do we need to encode it?

Obvious way

4 chars $\Rightarrow$ 2 bits/char $\Rightarrow$ 2n bits
                                        for  n chars

Can we do better?

Need more information

# Data Compression Problem -2

- Suppose we also know the frequency with which each character appears:

char:    A      B      C      D

freq:    .4     .3     .2     .1

- Can we use shorter codes (fewer bits) for more common characters?

Try $A=0$, $B=1$, $C=00$, $D=01$

what is $000$? $AAA$? $AC$? $CA$?

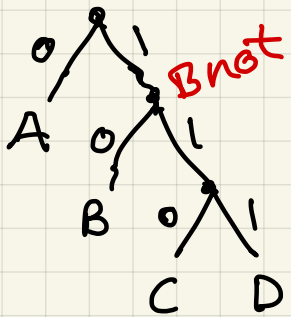Goal: avoid common "prefix", $0$ for A and C

3

# Prefix-free Codes

| char | A | B | C | D |
|------|---|---|---|---|
| freq | .4 | .3 | .2 | .1 |
| code | 0 | 10 | 110 | 111 |

total # bits for $n$ chars

$$= n \cdot .4 \cdot 1 + n \cdot .3 \cdot 2 + n \cdot .2 \cdot 3 + n \cdot .1 \cdot 3$$

$$= 1.9n \quad \text{versus } 2n \text{ for } 2 \text{ bit encodings}$$

Def: Code is prefix free if no char is prefix of another

Fact: 1-1 correspondence between prefix free codes and full binary trees (i.e. nodes have 0 or 2 children) with characters at leaves

4

# Prefix free Codes -2
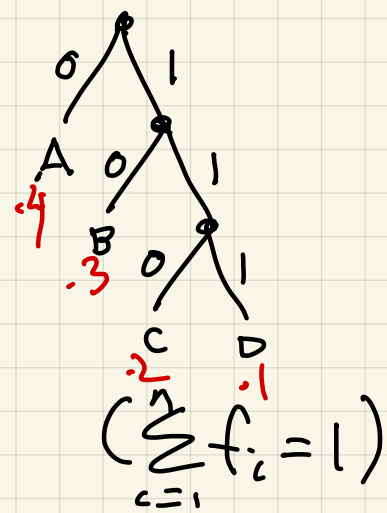
Goal: Find full binary tree with minimal cost, given

$n = \#$chars, $f_i = $ frequency of char $i$   $\left(\sum_{i=1}^{n} f_i = 1\right)$

$$\text{cost} = \sum_{i=1}^{n} f_i \cdot (\text{length of encoding of char } i)$$

$$= \sum_{i=1}^{n} f_i \,(\text{depth of leaf } i \text{ in tree})$$

sending $m$ chars requires $m \cdot \text{cost}$ bits

Greedy Intuition: Who goes at bottom of tree?
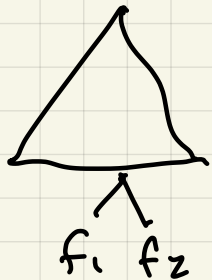i.e. longest encodings? least frequent chars.
smallest $f_i$ at bottom

5

# Cost of an Optimal Tree

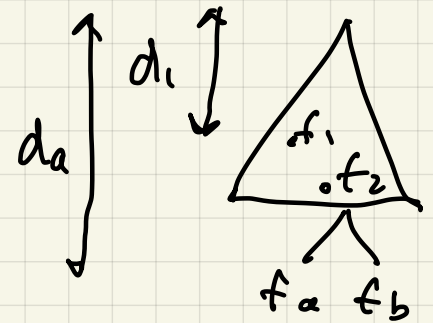$$Cost = \sum_{i=1}^{n} f_i \cdot (depth \ of \ leaf \ i \ in \ tree)$$

Claim: Suppose $f_1 \leq f_2 \leq \cdots \leq f_n$

$f_1$ and $f_2$ are siblings at bottom of tree.

proof: suppose not

then swap $f_1$ & $f_a$
$f_2$ & $f_b$

$d_a \quad d_1$

$f_1$
$f_2$

$f_a \quad f_b$

Claim lowers cost! swap $f_1$ & $f_a$ changes cost by

cost $=$
$3 \cdot f_1 + 3 \cdot f_2 + 2 \cdot f_3 + f_4$

$$(d_1 \cdot f_1 + d_a \cdot f_a) - (d_1 f_a + d_a f_1) \quad \text{apply}$$
$$= (d_a - d_1) \cdot (f_a - f_1) \geq 0 \quad \text{repeatedly?}$$

$\geq 0 \qquad \geq 0$

$f_4 \quad f_1 + f_2 + f_3$
$f_3 \quad f_1 + f_2$
$f_1 \ f_2$
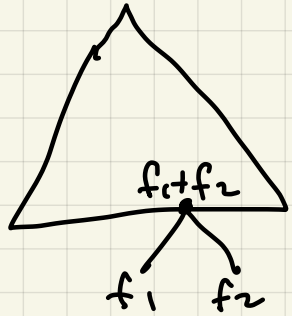
Def for any node except root, $cost(v) = \sum f_i$ of all leaves in subtree rooted at $v$

Claim Cost of tree $= \sum_{v \ except \ root} cost(v)$

6

# Greedy Algorithm for an Optimal Tree -1

1. pick 2 least frequent chars $f_1 \leq f_2 \leq \ldots$
2. replace them by 1 char with frequency $f_1 + f_2$
3. find optimal tree with $n-1$ chars
4. replace leaf for $f_1 + f_2$ by $\bigwedge_{f_1 \ f_2}$

proof: Induction on $n = \#$chars
    Base: $n = 1$ or $2$ chars, use $1$ bit

$T' = $ tree found on $f_1 + f_2, f_3, f_4 \ldots f_n$ $\left(\begin{array}{c}\text{optimal by} \\ \text{induction}\end{array}\right)$
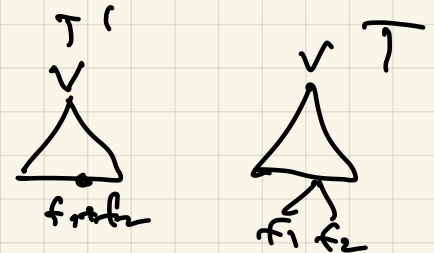
$T = T'$ with $\bigwedge_{f_1 \ f_2}$ added to $f_1 + f_2$

if $v \in T$ and $T'$ then $cost(v)$ same
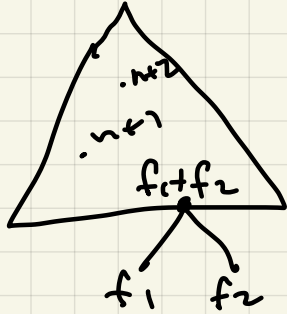
$\Rightarrow cost(T) = cost(T') + f_1 + f_2$

if $T$ not optimal $\Rightarrow \exists$ better tree for $f_1, f_2, f_3 \ldots$

$\Rightarrow T'$ not optimal (replace $f_1, f_2$ by $\bigwedge_{f_1 \ f_2}$    siblings at bottom

                       contradiction $\lightning$

# Greedy Algorithm for an Optimal Tree -2

1. pick 2 least frequent chars $f_1 \leq f_2 \leq \ldots$
2. replace them by 1 char with frequency $f_1 + f_2$
3. find optimal tree with $n-1$ chars
4. replace leaf for $f_1 + f_2$ by $\wedge$ $f_1$ $f_2$

Create Priority Queue $H$ of $\{1, 2, \ldots, n\}$
ordered by $f_1 \leq f_2 \leq \ldots \leq f_n$

for $k = n+1$ to $2n-1$

$\begin{cases} i = \text{deletemin}(H), \quad j = \text{deletemin}(H) \\ \text{create node } k \text{ with children } i \text{ and } j \\ f_k = f_i + f_j \; ; \; \text{insert}(k) \end{cases}$

binary heap

cost of A $lg = O(n \cdot \text{cost}(\text{deletemin}) + n \cdot \text{cost}(\text{insert})) = O(n \log n)$
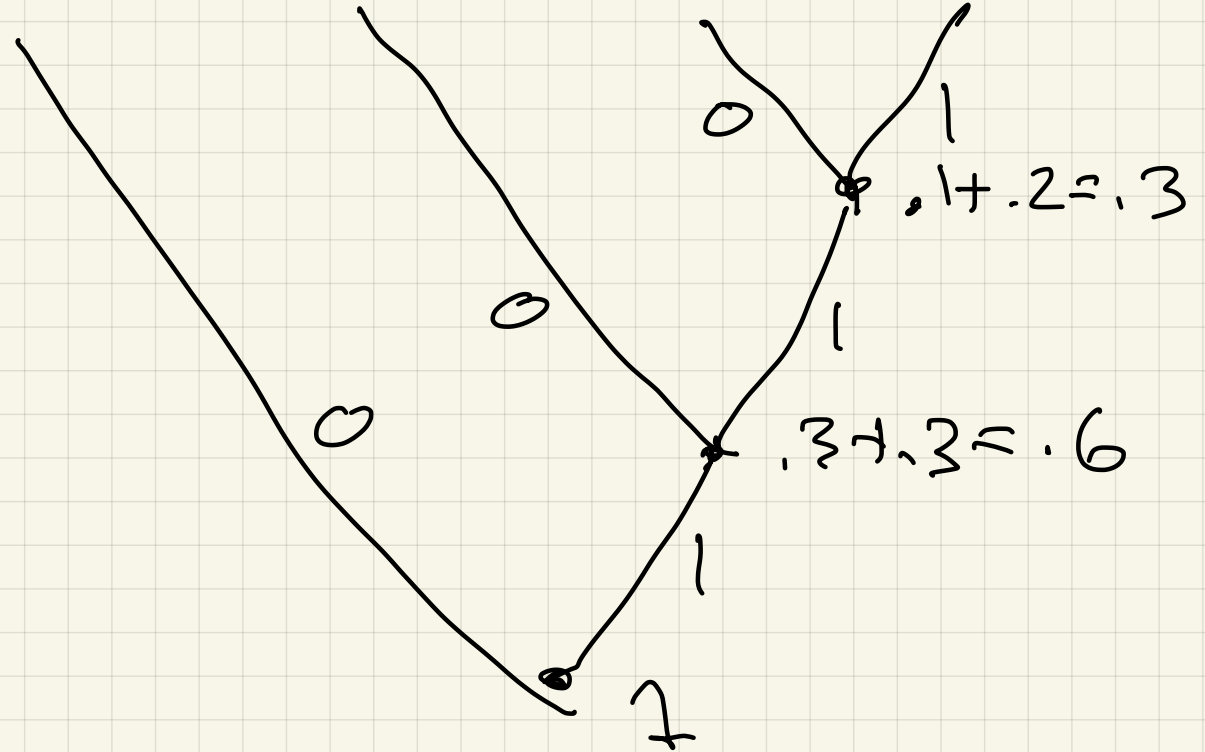
8

# Huffman Encoding Example

| char | A | B | C | D |
|------|---|---|---|---|
|      | 0 | 10 | 110 | 111 |
| freq | .4 | .3 | .2 | .1 |



$.1 + .2 = .3$

$.3 + .3 = .6$

# Horn Formulas

- Can we satisfy a Boolean expression?

  - Notation $w$ = Boolean variable (T or F)
  
    $\wedge$ = and, $\vee$ = or, $\bar{w}$ = not w

  - Are there values (T or F) of Bool. vars. that make this True?

    $$(w \wedge y \wedge z) \wedge (x \wedge y \Rightarrow w) \wedge \underbrace{(\bar{u} \vee \bar{v} \vee \bar{z})}_{clause} \cdots$$

- General case: NP Hard (chap 8)

- Horn Formulas — special case — clauses must b

  1) implication $(u \wedge v \wedge w) \Rightarrow x$

     with all positive vars (no $\bar{u}, \bar{v}$ etc)

     includes $\Rightarrow x$ ($x$ most be T)

     recall $a \Rightarrow b \equiv \bar{a} \vee b$

  2) negative clauses $(\bar{x} \vee \bar{y} \vee \bar{z})$

10

# Horn Formula : Example

$x \equiv$ murder took place in kitchen

$y \equiv$ butler innocent

$z \equiv$ colonel asleep at 8pm

$w \equiv$ murder took place at 8pm

$u \equiv$ colonel innocent

$v \equiv$ professor innocent

implication $(z \wedge w) \Rightarrow u$

negative clause $(\bar{u} \vee \bar{y} \vee \bar{v})$

See Prolog

# Greedy Algorithm for Horn Formulas

3 kinds of formulas:

1) $(z \wedge w) \Rightarrow u$

2) $\qquad\qquad \Rightarrow x \quad (x = T)$

3) $(\bar{u} \vee \bar{v} \vee \bar{y})$

Set all variables to false $\left( \begin{array}{l} 3) \text{ ok} \quad 2) \text{ not ok} \\ 1) \text{ ok} \end{array} \right)$

While an implication not satisfied, set rhs = T

If all negative clauses satisfied

    Satisfiable: return assignment

else

    Not satisfiable

Example:

$(w \wedge y \wedge z) \Rightarrow x, (x \wedge z) \Rightarrow w, x \Rightarrow y, \Rightarrow x, (x \wedge y) \Rightarrow w, (\bar{w} \vee \bar{x} \vee \bar{y})$

    OK          OK       T   T          T     oops  12

# Back to Huffman: Intro to Entropy

- What is Entropy?

- What is "Information"?

  - If a random event occurs, and I tell you the outcome, how much Information is that?

  - $I(p)$ = information from being told a random event with probability $p$ occurred

- Base case: flip fair coin, be told H or T : $I(\frac{1}{2}) = 1$ bit

- If $p > \frac{1}{2}$, $I(p) < 1$, if $p < \frac{1}{2}$, $I(p) > 1$

- Two indep. events, probs $p_1$ and $p_2 \Rightarrow I(p_1 \cdot p_2) = I(p_1) + I(p_2)$

$$\Rightarrow I(p) = \log_2 \left( \frac{1}{p} \right)$$

# Intro to Entropy -2

- Suppose $n$ possible outcomes, probabilites $p_1, p_2, \ldots, p_n$

$$\text{Expected information} = E(I) = \sum_{i=1}^{n} p_i I(p_i) = \sum_{c=1}^{n} p_i \log_2\left(\frac{1}{p_i}\right)$$

$$= \text{Entropy, measures how "random" distribction is}$$

Ranges from $0$ (one $p_i = 1$, rest $0$) to $\log_2 n$ (all $p_i = \frac{1}{n}$)

- In Huffman, suppose all $f_i = p_i = \frac{1}{2^{k_i}}$ for some $k_i$
  - can show depth of $f_i$ in Huffman tree $= k_i$
  - To encode $m$ chars with frequencies $p_i$, takes
  
  $$m \cdot \sum_{i=1}^{n} p_i \cdot \text{depth of } f_i \text{ in tree} = m \cdot E(I) = m \cdot \text{Entropy}$$

- Thm (Shannon) $m \cdot$ Entropy is a lower bound on any encoding scheme

14