

CS 223 Lec 16

Question: how do we find
an optimal prefix code?

Answer: Huffman trees

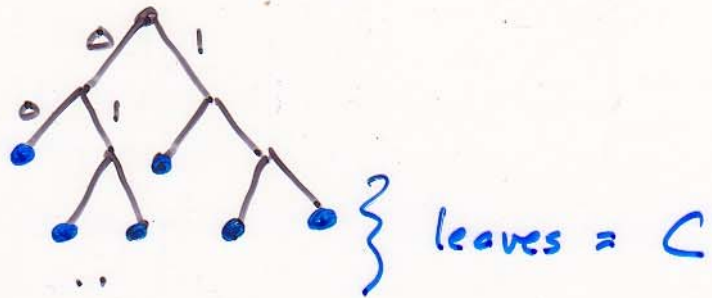
Let F be a text file,
over an alphabet C .

It is easy to compute
the frequency of each
character $c \in C$.

$f[c] = \#$ of occurrences
of c in F

Let T be a tree
representation of a prefix
code.

T



If we use T to encode F
 how long is the
 bit string produced?

$$B(T) = \# \text{ of bits used}$$

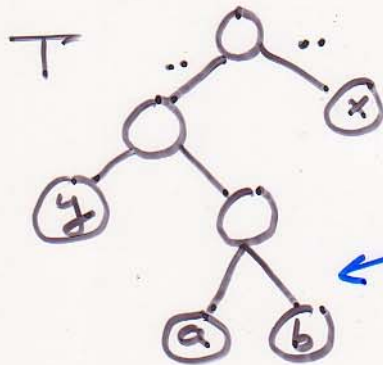
$$= \sum_{c \in C} f(c) \cdot \text{depth}_T(c)$$

goal: find a tree T
 that minimizes $B(T)$

Lemma Let $x, y \in C$ that have lowest frequencies.

\exists an optimal tree T where x and y are siblings.

Proof Let T be some optimal tree. Suppose x and y are not currently siblings.

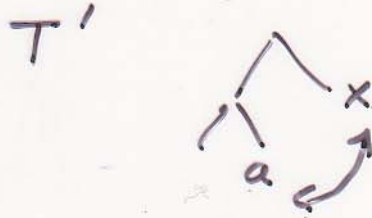


\leftarrow a and b are two nodes that are siblings and at maximum depth.

Observe we can assume

$$f(a) \leq f(b) \quad \text{and} \quad f(x) \leq f(y)$$

Let T' be the tree obtained
by exchanging a and x



$$\begin{aligned} B(T) - B(T') &= \sum_c f(c) d_T(c) - \sum_c f(c) d_{T'}(c) \\ &= f(x) d_T(x) + f(a) d_T(a) \\ &\quad - f(x) d_{T'}(x) - f(a) d_{T'}(a) \\ &= f(x) (d_T(x) - d_{T'}(x)) \\ &\quad + f(a) (d_T(a) - d_{T'}(a)) \\ &= \cancel{f(x) (d_T(x) - d_T(a))} + \cancel{f(a) (d_T(a) - d_T(x))} \\ &= f(x) (d_T(x) - d_T(a)) \\ &\quad + f(a) (d_T(a) - d_T(x)) \end{aligned}$$

$$= \underbrace{(f(x) - f(a))}_{\leq 0} \underbrace{(d_T(x) - d_T(a))}_{\leq 0}$$

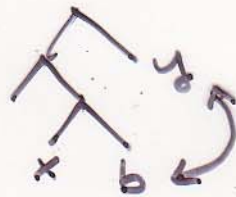
$$\geq 0$$

$$\Rightarrow B(T) \geq B(T')$$

↑
optimal

So, T' must be optimal.

T'' :



It is easy to show

$$B(T'') \leq B(T')$$

So T'' is also optimal.

Idea: merge the

two least frequently occurring characters (nodes)

to form a new node

whose freq. is the sum

of two children.

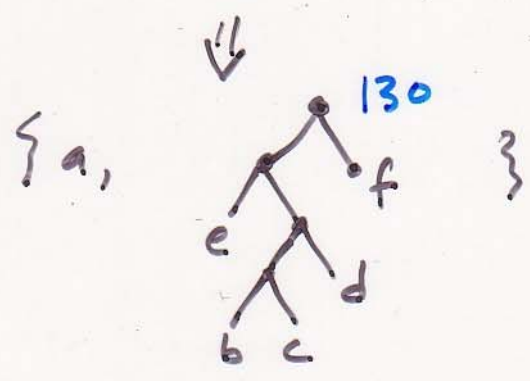
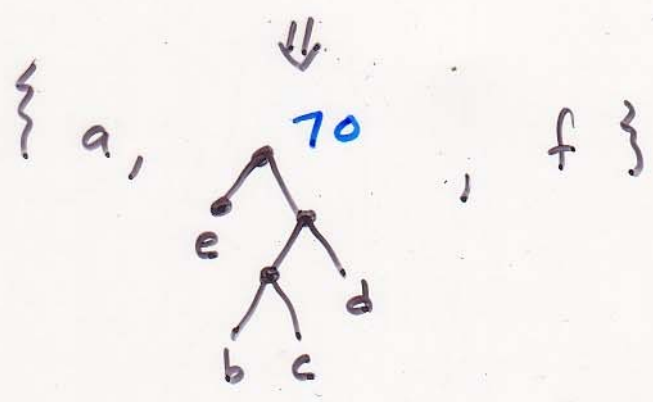
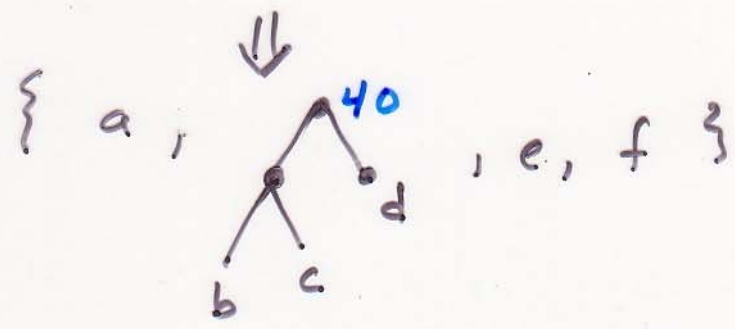
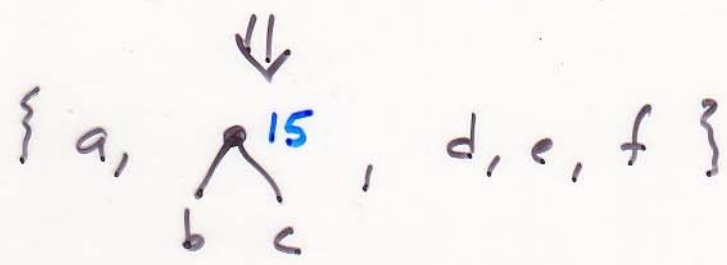
repeat

Example

$$C = \{ a, b, c, d, e, f \}$$

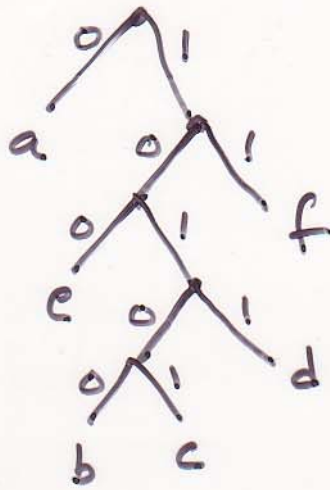
	f
a	100
b	10
c	5
d	25
e	30
f	60

{ a, b, c, d, e, f }





optimal tree \Rightarrow



<u>freq</u>	<u>code</u>
100	a 0
10	b 10100
5	c 10101
25	d 1011
30	e 100
60	f 11

$$\begin{aligned}
 B(T) &= \# \text{ of bits used} \\
 &= 100 \cdot 1 + 10 \cdot 5 + 5 \cdot 5 + 25 \cdot 4 + \\
 &\quad 30 \cdot 3 + 60 \cdot 2 \\
 &= 485 \text{ bits.}
 \end{aligned}$$

Lemma Huffman's procedure leads to an optimal code.

Pf Form a tree by recursively merging the two least freq. nodes ... results in optimal tree



T be the optimal tree for original problem

T' the optimal tree for alphabet $C = \{x, y\} + \{z\}$

$\Rightarrow T$ is optimal provided

T' is \Rightarrow establish that Huffman's procedure is correct.