# CS 314 - Huffman codes & shortest paths

## Announcements

# Huffman codes

Goal: Transmit a message using as few bits as possible.

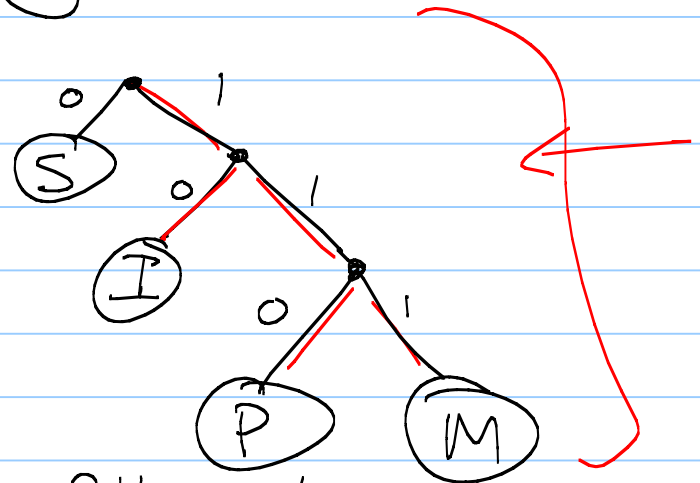[ Use frequency counts (so know the message ahead of time).

Huffman codes: <u>pre-fix free</u>

Why? So we can scan & decode — no ambiguity.

Can visualize as a binary tree

M: 111
I: 10
S: 0
P: 110



When reading string, just follow tree & output letter when you reach a leaf.

Ex: 0 10 11111 10
SIMP

Given n letters plus frequency counts
for each letter, $f[1..n]$ ←
find the code minimizing total length:

$$\sum_{i=1}^{n} f[i] \cdot depth(i) \quad = cost(T)$$

# Pseudocode:

→ Keep characters in min heap, w/ priority = to frequency.

[ L, R, & P keep track of left/right and parent indices.

```
BUILDHUFFMAN(f[1..n]):
    for i ← 1 to n
        L[i] ← 0;  R[i] ← 0
        INSERT(i, f[i])

    for i ← n to 2n − 1
        x ← EXTRACTMIN()
        y ← EXTRACTMIN()
        f[i] ← f[x] + f[y]
        L[i] ← x;  R[i] ← y
        P[x] ← i;  P[y] ← i
        INSERT(i, f[i])

    P[2n − 1] ← 0
```

$O(n \lg n)$

] get two least frequent letters

← sum frequencies

← make them leaves

## Proof of Correctness:

Lemma: Let $x$ & $y$ be the 2 least frequent characters. Then there is an optimal code where $x$ & $y$ are siblings and have maximum depth in the tree.

Did proof in class last time.

Thm: Huffman codes are optimal.

Pf: Induction on # of letters.

Base case: $n = 1$ (or 2) ✓

IH: Given $< n$ characters, Huffman's alg. is optimal.

IS: $n$ characters with frequency counts $f[1..n]$ wlog, assume $f[1]$ & $f[2]$ are least frequent.
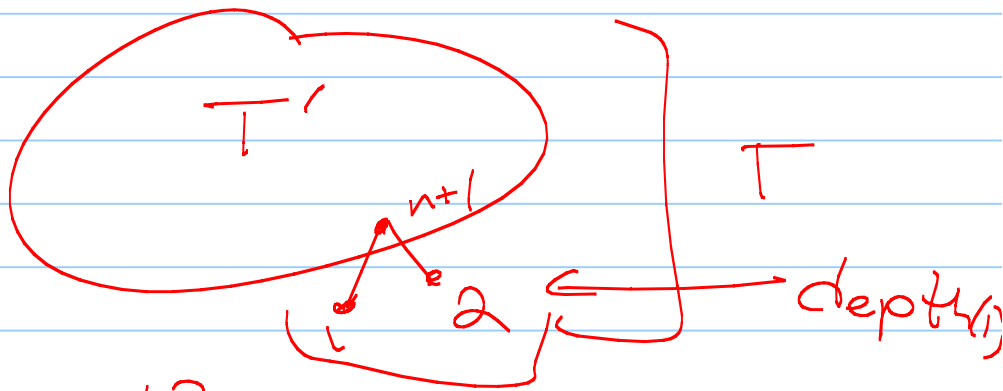By our prev. lemma, some optimal tree has $1$ & $2$ as siblings (at deepest level).

pf continued : Create $f[n+1] = f[1] + f[2]$.
    Let $T'$ be Huffman tree for $f[3..n+1]$.
→ By IH, $cost(T')$ is smallest possible
    for any such tree.
Create $T$ by removing leaf $n+1$ &
replacing it w/ internal node w/ 2 children,
$1 + 2$.



Why is $T$ optimal?

**Claim:** T is optimal for $f[1..n]$.

$$\text{cost}(T) = \sum_{i=1}^{n} f[i] \cdot \text{depth}(i)$$

$$= \sum_{i=3}^{n+1} f[i] \cdot \text{depth}(i) - f[n+1] \cdot \text{depth}(n+1) + f[1] \cdot \text{depth}(1) + f[2] \cdot \text{depth}[2]$$

$$= \text{cost}(T') + (f[1] + f[2]) \cdot \text{depth}(1)$$
$$- f[n+1] \cdot \text{depth}(n+1)$$

$$(\text{know } f[1] + f[2] = f[n+1], \text{ and}$$
$$\text{depth}(1) = \text{depth}(n+1) + 1)$$

$$= \text{cost}(T') + f[1] + f[2].$$

So $\boxed{\text{cost}(T) = \text{cost}(T') + f[1] + f[2]}$.

Suppose $T$ was _not_ optimal.

In optimal tree
remove 1 & 2 & get a
Huffman tree for $3 \ldots \cup n+1$.

If do that, get a tree for $3 \ldots n+1$
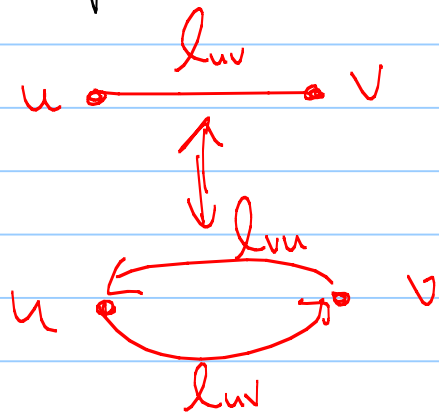that is ∪better than $T'$.

contradiction

# Shortest paths in a graph. (4.4)

Suppose we have $G = (V, E)$ and each edge $e \in E$ has a length $l_e$.

Here, we'll assume $G$ is directed: $u \bullet \longrightarrow \bullet v$.

If given undirected graph, how could we adapt to directed model?

$u \bullet \xrightarrow{\quad l_{uv} \quad} \bullet v$

$\updownarrow$

$u \bullet \underset{l_{uv}}{\overset{l_{vu}}{\rightleftarrows}} \bullet v$

$O(2|E|) = O(|E|)$

$s_0$ ——————————→•$t$

**Goal:** Given two vertices, find shortest path between them.

Why?    Mapquest!

$9 \diagup 1 \diagdown 2 \diagdown 25$

Ideas?



$L_1$
$L_2$
$j$

BFS — right idea

← distance 1

We'll actually do something harder:

Given a source vertex s, compute shortest path from s to every other vertex.

The reason — if we don't explore everything, we don't know if we've missed a shorter path.

# Greedy idea:

Start with a set S.
(initially $S = \{s\}$)

At each step, grow out from s,
taking next shortest path from s to
a new vertex + adding that to S.

Ex: