# CS2100:

## Huffman codes

# Recap

Lab due Sunday ←

Reading for Zybook —
check later/today (Graphs)

Next HW is posted

Final exam:

, Wed at 2pm

⌐→ Last day of class:
review session

# Huffman Codes - the idea:

We would like to transmit info using as few bits as possible.

What does <u>ASCII</u> do?

8-bit rep. of letters
$\hookrightarrow$ 256 characters

X letters $\Rightarrow$ 8x bits

How can we do better?
$\hookrightarrow$ Well, what if we don't <u>use</u> all the characters?

fewer than 8 bits
$\hookrightarrow$ shorter

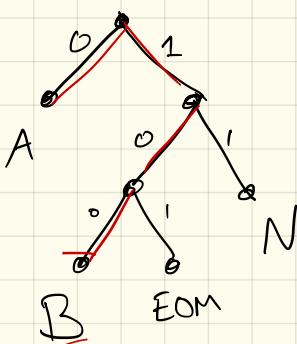use fewer bits for more common letters

# Problem:

If not fixed length, hard to tell when a character is finished.

Ex:
E : 11  ⎫
A : 00  ⎬ 2 bits
S : 01  ⎪
T : 10  ⎭
R : 110 ⎫
M : 001 ⎬ 3 bits
B : 010 ⎪
N : 100 ⎭
: etc.

Decode: 11001
        E  M
        R  S    ⎬ X

# Prefix-free codes

An unambiguous way
to send
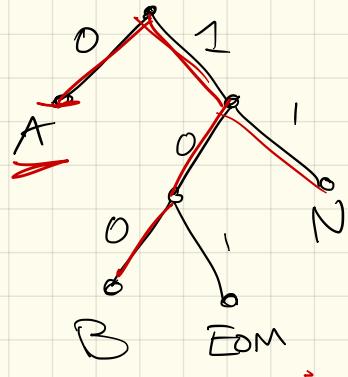information when
we have characters
not of a fixed
length.

Tree diagram with edges labeled:
- Root splits: 0 (to A) and 1
- From 1 node: 0 and 1 (to N)
- From 0 node: 0 (to B) and 1 (to EOM)

A

N

B    EOM

Key: No letter's code will
be the prefix of another.

Encode: BAN

100011

## Decode:

1 0 0 0 1 1 0 1 1 0 1 0 1

B A N A N A ↑ EOM



Bits : 13

ASCII: 6 × 8

# How should we do this?

Use frequency counts
to make a good
prefix-free code (or tree):

This sentence contains three a's, three c's, two d's, twenty-six e's, five f's, three g's, eight h's, thirteen i's, two l's, sixteen n's, nine o's, six r's, twenty-seven s's, twenty-two t's, two u's, five v's, eight w's, four x's, five y's, and only one z.

lower frequency
= more bits
↳ lower in tree

higher frequency
= fewer bits
↳ higher in tree

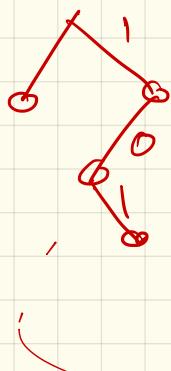# Goal: Minimize Cost

$\hookrightarrow$ here, minimize total
length of encoded
message:

Input: frequency counts
$$f[1..n]$$

Compute: tree with
minimum "cost"

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



\# of bits
to transmit

To do this, we'll need to
use the array f:

This sentence contains three a's, three c's, two d's, twenty-six e's, five f's, three
g's, eight h's, thirteen i's, two l's, sixteen n's, nine o's, six r's, twenty-seven s's,
twenty-two t's, two u's, five v's, eight w's, four x's, five y's, and only one z.

If we ignore punctuation
& spaces (just to keep
it simple),
we get:

| A | C | D | E | F | G | H | I | L | N | O | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 3 | 2 | 26 | 5 | 3 | 8 | 13 | 2 | 16 | 9 | 6 | 27 | 22 | 2 | 5 | 8 | 4 | 5 | 1 |

Which letters should
be deeper (or shallower)?
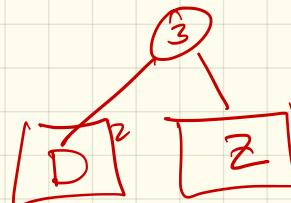
(ie: How to be greedy?)

# Huffman's alg:

Take the two least frequent characters.

Merge them into one letter, which becomes a new "leaf":

| A | C | D | E | F | G | H | I | L | N | O | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 3 | 2 | 26 | 5 | 3 | 8 | 13 | 2 | 16 | 9 | 6 | 27 | 22 | 2 | 5 | 8 | 4 | 5 | 1 |

| A | C | E | F | G | H | I | L | N | O | R | S | T | U | V | W | X | Y | DZ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 3 | 26 | 5 | 3 | 8 | 13 | 2 | 16 | 9 | 6 | 27 | 22 | 2 | 5 | 8 | 4 | 5 | 3 |

# Example (cont):

| AC / 6 | A | C | E | F | G | H | I | | N | O | R | S | T | | | W | X | Y | Z | LU / 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 3 | 3 | 26 | 5 | 3 | 8 | 13 | 2 | 16 | 9 | 6 | 27 | 22 | 2 | 5 | 8 | 4 | 5 | 6/8 | |

# In the end, get a tree with letters at the leaves:



A Huffman code for Lee Sallows' self-descriptive sentence; the numbers are frequencies for merged characters

| A | C | D | E | F | G | H | I | L | N | O | R | S | T | U | V | W | X | Y | Z |
|---|---|---|----|---|---|---|----|---|----|---|---|----|----|---|---|---|---|---|---|
| 3 | 3 | 2 | 26 | 5 | 3 | 8 | 13 | 2 | 16 | 9 | 6 | 27 | 22 | 2 | 5 | 8 | 4 | 5 | 1 |

If we use this code, the encoded message starts like this:

1001 0100 1101 00 00 111 011 1001 111 011 110001 111 110001 10001 011 1001 110000 ...

 T    H    I   S  S  E   N   T    E   N   C      E   C      O     N   T    A

# How many bits?

| char. | A | C | D | E | F | G | H | I | L | N | O | R | S | T | U | V | W | X | Y | Z |
|-------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| freq. | 3 | 3 | 2 | 26 | 5 | 3 | 8 | 13 | 2 | 16 | 9 | 6 | 27 | 22 | 2 | 5 | 8 | 4 | 5 | 1 |
| depth | 6 | 6 | 7 | 3 | 5 | 6 | 4 | 4 | 7 | 3 | 4 | 4 | 2 | 4 | 7 | 5 | 4 | 6 | 5 | 7 |
| total | 18 | 18 | 14 | 78 | 25 | 18 | 32 | 52 | 14 | 48 | 36 | 24 | 54 | 88 | 14 | 25 | 32 | 24 | 25 | 7 |

Total is  $\sum f[i] \cdot depth(i)$

$= 646$  bits  here

How would ASCII do on these
170  letters

$\hookrightarrow 170 \times 8$

**Thm:** Huffman codes are optimal:
they use the fewest # of bits
possible.

pf:   (go take 3100)

Side note: This is known as
a greedy algorithm.

# Implementation: use priority queue

**heap!**

---

**BUILDHUFFMAN($f[1..n]$):**
  for $i \leftarrow 1$ to $n$
      $L[i] \leftarrow 0$; $R[i] \leftarrow 0$
      INSERT($i, f[i]$)

  for $i \leftarrow n$ to $2n-1$
      $x \leftarrow$ EXTRACTMIN( )
      $y \leftarrow$ EXTRACTMIN( )
      $f[i] \leftarrow f[x] + f[y]$
      $L[i] \leftarrow x$; $R[i] \leftarrow y$
      $P[x] \leftarrow i$; $P[y] \leftarrow i$
      INSERT($i, f[i]$)

  $P[2n-1] \leftarrow 0$

---

+ 3 arrays  $L, R, P$:

(go take algorithms)

# Next HW:

## decode:

Given an input which describes a tree & a message:

1) Create the tree
2) Use it to decode the message.

One thing I skipped: do need to store the tree.

Overview of assignment...