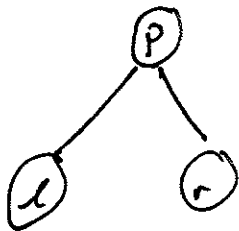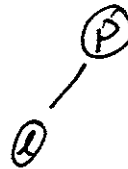# Binary Trees

- A tree is a graph $T = (V, E)$ that is acyclic (contains no cycles)

  - between any 2 vertices there is exactly __1__ paths
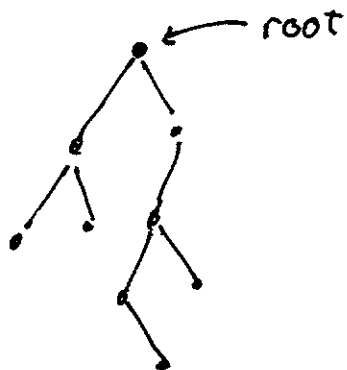
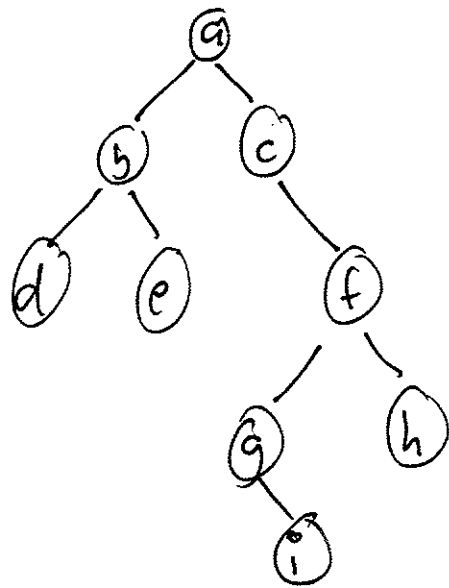  - binary tree: oriented, parent, left, right child

no children: leaf

- a vertex with no parent

root

root to any node: 1 path

~~ℓ~~ length of that path $r \rightsquigarrow u$ is the depth of $u$

| node | depth | height |
| --- | --- | --- |
| f | 2 | 2 |
| c | 1 | 3 |
| a | 0 | 4 |
| i | 4 | 0 |

Depth of T : max. depth of any node
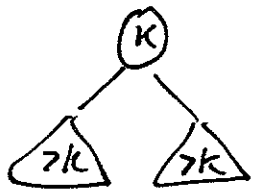
Depth : 4

Height of a node u is the longest path to
   any of its descendent leaves

# Binary Search Tree:



# Heap: (min heap)

- Binary tree

- K in a node $u$ is less than both its children



- Full: every node is present except the last level, full to the left

# Coding Theory:

- encoding symbols (in binary)

- ASCII: fixed length encoding   all codewords have length 8

$$A = 65 = \underbrace{01000001}_{8 \text{ bits}}$$

- Variable length encoding:
  - common letters: shorter code words $\Big\}$ reduces the average code word length
  - rare letters: ~~~~ longer code words
  - omit some characters

$$\begin{matrix} A & B & C & D \end{matrix}$$

$\{0, 01, 101, 010\}$    variable length code

message: $010 \longrightarrow D$    ambiguous

$\longrightarrow BA$

Goal: a <u>prefix free code</u>: no code word is the prefix

of another

A is a prefix of B   0 vs 01

B is a prefix D

**Idea:**

- Want to encode a fix usig a variable lengh, prefix tree code That minimizes th aveye code word length (compresion) loss less

- Huffman Coding

  - high frequeny (common) letters: small code words
  - low frequy (rare) letters: lonyer code words

- Given g frequencies, build a tree:

  combine "small" frequenya as children to a new root node ⟶ new tree

- continue until 1 tree remains

- path root ⟶ leaf defines a code word.

Fixed length encoding: 3

| Symbol | codeword | length + freq |
|--------|----------|---------------|
| A | 101 | 3 , .10 |
| B | 001 | 3 , .15 |
| C | 01 | 2 , .34 |
| D | 1001 | 4 , .05 |
| E | 000 | 3 , .12 |
| F | 11 | 2 , .21 |
| G | 1000 | 4 , .03 |

```
0 0 0
0 0 1
0 1 0
0 1 1
1 0 0
1 0 1
1 1 0
TH
```

$$\frac{3 - 2.53}{3} = 15.67\%$$
(compression ratio)

Average code word length

pack

Zip $\neq$ Huffman

$.3 + .45 + .68 + .2 + .36 + .42 + .12$

$= 2.53$

Input: Alphabet of symbols, $\Sigma$

$\quad$ frequency of symbols $freq: \Sigma \rightarrow \mathbb{R}^+$ $\qquad$ $freq(x) =$ frequency of the

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ symbol $x$ in the file

Output: A Huffman Tree

$H \longleftarrow$ min Heap

//initialization

for each $x \in \Sigma$

$\quad T_x \longleftarrow$ single node tree

$\quad wt(T_x) \longleftarrow freq(x)$

$\quad$ H.insert $(T_x)$

while $H.size > 1$ $\qquad$ $(n-1)$

$\quad T_p \longleftarrow$ new root

$\quad T_\ell \longleftarrow$ H.getMin

$\quad T_r \longleftarrow$ H.getMin

$\quad T_p.leftChild \longleftarrow T_\ell$

$\quad T_p.rightChild \longleftarrow T_r$

$\quad wt(T_p) \longleftarrow wt(T_\ell) + wt(T_r)$

$\quad$ H.insert $(T_p)$

Output H.getMin

getMin, insert:

$\quad O(d) = O(\log(n))$

$n =$ number of things
$\qquad$ in the heap

$$\log(a) + \log(b)$$
$$= \log(a \cdot b)$$

total cost:

$$\sum_{k=2}^{n} \log(k) = \log(2) + \log(3) + \log(4) + \cdots \log(n)$$

$$= \log(2 \cdot 3 \cdot 4 \cdots n)$$

$$= \log(n!)$$

$$\leq \log(n^n)$$

$$= n \log(n)$$

$$\Theta(n \log(n))$$

| iteration | size of H | getMin/insert cost |
|---|---|---|
| 1 | n | $\log(n)$ |
| 2 | $n-1$ | $\log(n-1)$ |
| 3 | $n-2$ | $\log(n-2)$ |
| $\vdots$ | | |
| $i$ | $n-i(+1)$ | $\log(n-i)$ |
| $\vdots$ | | |
| ~~n-1~~ $n-1$ | 2 | $\log(2)$ |