# Greedy strategies
## (Huffman codes)

→ Encode $\{a, \dots, z\}$ using $\{0, 1\}$

will need 5 bits (32 5 bit strings)
if each letter is a 5 bit string

→ Can we optimize amount of data to transfer?

If $n$ letters $\Rightarrow$ $5n$ bits

→ Variable length encoding

(Morse code)          •          —

          o          1

freq      $\begin{cases} e - 0 \\ t - 1 \end{cases}$
letters

next      $\{a - 01$
freq

But decoding ??

$0101 \rightarrow aa$

$\swarrow$  $\searrow$

et et                          et a

$\cdot\text{-}\text{-}$

→ In practice , " slight pause"          (ambiguous)

between    letters.

→    so    morse code is          $\cdot$, $\text{---}$,    and , pause

(3   alphabet)

---

want    Variable  length  coding   with   Prefix  code property

* E(x)    =    encoding ( x )

* E(x)    not    a    prefix for any   E(y) ]

prefix code
property

(eg)                x | a | b | c | d | e

$$F(x) \mid 11 \mid 01 \mid 001 \mid 10 \mid 000$$

$$\underline{001} \quad \underline{000} \quad \underline{001} \quad \underline{11} \quad \underline{01}$$
$$\phantom{00}c \quad\quad e \quad\quad c \quad\quad a \quad b$$

(unambiguous decoding possible)

---

## Optimal prefix codes

Given $\quad f(x) = $ freq of $x$    (statistical analysis)

$$\sum_{x \in \text{alphabet}} f(x) = 1 \qquad \forall x$$

$\rightarrow$ n letter message

$\rightarrow$ so $\approx$ $nf(x)$ # of occurences $\forall x$

$\rightarrow \qquad \sum_{\substack{x \in \\ \text{alphabet}}} \text{Length}[E(x)] \cdot n \, f(x) = $ # of bits to encode n letter message.

(Expected #)

→ Avg # of bits / letter

$$= \sum_{\substack{x \in \\ \text{alphabet}}} |E(x)| f(x)$$

$$\left[ \begin{array}{l} \text{Fixed length code:} \quad \text{each letter is encoded by} \\ \quad m \quad \text{bits} \quad \forall \quad \text{letters in alphabet,} \quad |E(x)| = m \end{array} \right]$$
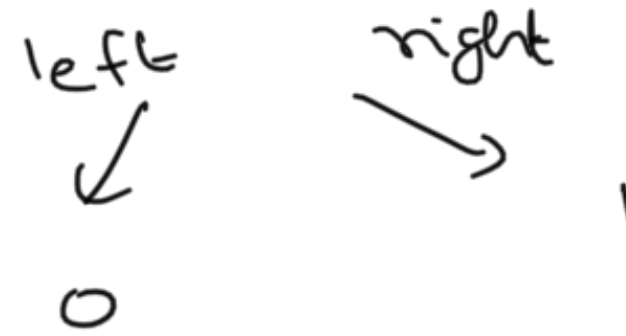
Idea: shorter codes for more frequent letters....

Find $E(x)$ ∋: ① $\sum_{\substack{x \in \\ \text{alphabet}}} |E(x)| f(x)$ is minimized
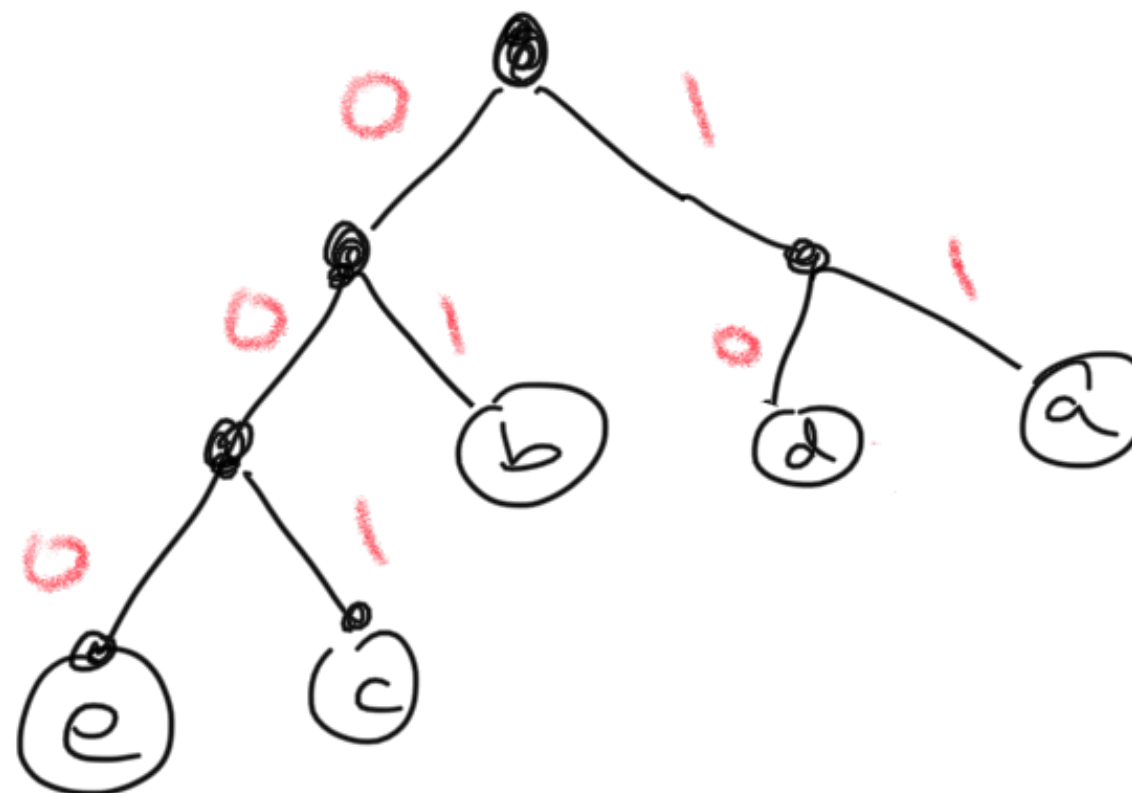
② has prefix code property

# Codes as trees

binary tree:     left     right

0    1

path from root to a leaf: binary seq

a    b    c    d    e

11    01    001    10    000

⟱



Prefix code ⟹ each letter is at a leaf

property

$$\text{depth}(`c') = \text{length}(001)$$
$$= 3 = |E(`c')|$$

So want to put higher freq letters at lower depths
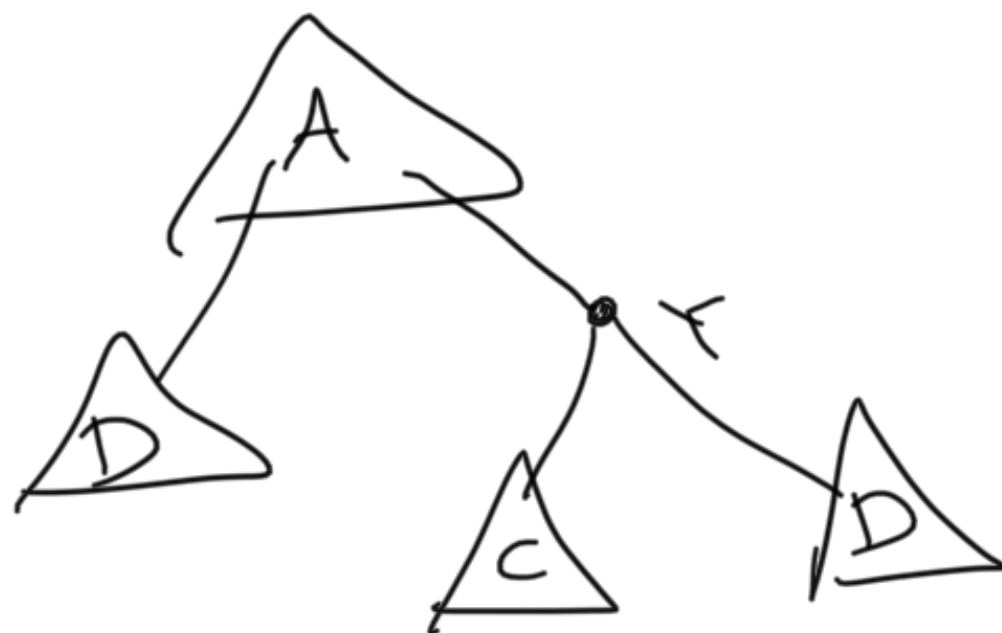
Full tree: 0 or 2 children for each node

I Optimal prefix code $\Rightarrow$ its tree has to be full

Pf



← if only child

shrink tree...

II     optimal prefix code $\Rightarrow$     If   depth$(x) \leq$ depth$(y)$
                                             then   $f(x) \geq f(y)$
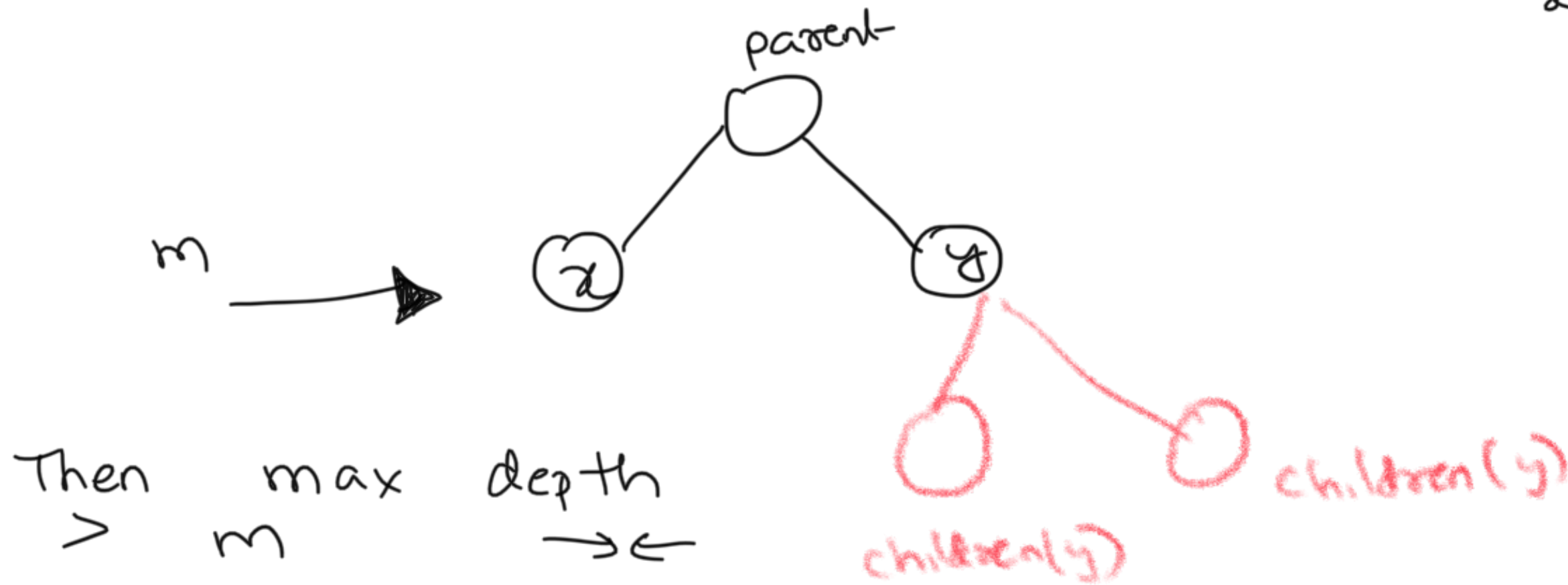
Pf: Else   exchange   x   and   y

III     Let   max depth of tree = m
        from   optimal prefix code

Let $x$ be a leaf @ depth $m$

Its sibling is also a leaf

Pf

since full,
$x$ has
sibling $y$

parent

$m \longrightarrow$ $x$     $y$

children($y$)

children($y$)

Then max depth
> $m$          $\rightarrow\leftarrow$

So Leaves @ max depths occurs in pairs

So leaves @ max depth $\rightarrow$ are lowest

# Recursion  (Huffman coding)

→ choose 2 lowest freq   ... $x, y$ ...

sibling leaves @ max depth



?? 

Alphabet  $A$

← treat as a unit

create new alphabet letter

change alphabet

$A \longrightarrow A' = A \smallsetminus \{x, y\} \cup \{\, \boxed{xy} \,\}$
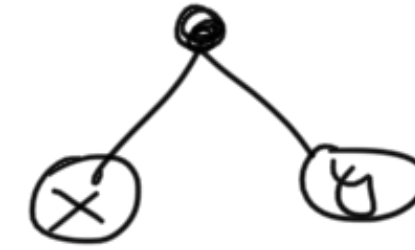
$f(\, \boxed{xy} \,) = f(x) + f(y)$

→ So recurse on  $A'$

→ Base case :  $|A| = 2$   $\Rightarrow E(x) = 0$

$$A = \{x, y\}$$

$$E(y) = 1$$

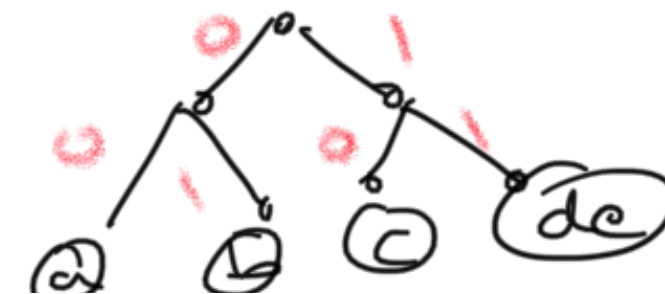then  go  from  $A^l$  $\longrightarrow$  A  by

replacing  (xy)  with

example

| x | a | b | c | d | e |
|---|---|---|---|---|---|
| f(x) | 0.32 | 0.25 | 0.2 | 0.18 | 0.05 |

$\longrightarrow$

| x | a | b | c | de |
|---|---|---|---|---|
| f | 0.32 | 0.25 | 0.2 | 0.23 |

| x | a | b | c de |
|---|---|---|---|
| f | 0.32 | 0.25 | 0.43 |

| x | a b | c de |
|---|---|---|
| f | 0.57 | 0.43 |

## why algo work?

→ Base case optimal

→ Assume ok for $k-1$ letters $|A'|$

→

$$A, T \quad \rightsquigarrow \quad A', T'$$

$$x, y \quad \cdots \qquad \boxed{xy} \quad \cdots$$

$$E = \sum f(char) |E(char)|$$

alphabet    tree
$A'$         $T'$

$$E' = \sum f(char') \, |E(char')|$$

$$E - E' = f(x)|E(x)|$$
$$+ f(y)|E(y)|$$
$$- f(xy)|E(xy)|$$

$$= f(x)|E(x)| + f(y)|E(y)|$$
$$- [f(x) + f(y)]|E(xy)|$$

$$= f(x)\ell + f(y)\ell$$
$$- f(x)(\ell-1) - f(y)(\ell-1)$$

$$= f(x) + f(y)$$

$$|E(x)| = \ell$$
$$|E(y)| = \ell$$
$$|E(xy)| = \ell-1$$

$$= f(xy)$$

* Suppose $\exists$ S (better tree) with $E_S > E_T$

\*     $x, y$     with     lowest    $f(x), f(y)$

     must   be   @   max   depth   in   S   also

\*     wlog     assume      (figure: tree with $x$ and $y$ as siblings)    in   S    also

     ( if   not , move   labels   @   max   depth

     leaves    to    make   x   and   y    as

     siblings )

\*     merge    $x, y$   $\longrightarrow$    $xy$   ,    $f(xy) = f(x) + f(y)$

     to    set   $S'$     ( / alphabet   $A'$

                                     of   $k-1$   letters )

$S' / A'$                            $T' / A'$                 $|A'| = k-1$

                                                 optimal

$$E_{S'} \geq E_{T'}$$

$$E_S = E_{S'} + f(xy) = E_{T'} + f(xy) = E_T$$

so T optimal as well

## Implementation

* Extract lowest 2 freq

* merge, replace by combined freq

Each recursive step

current
↓
* takes $O(|A|)$

# of recursive calls is $(|A| = k) - 1$

$$= k-1$$

* Store freq in array
* Scan to find min.1, min 2

$$2 |A|$$
$$= O(|A|)$$

$$k-1 + k-2 + \cdots + 1 = \frac{k(k-1)}{2}$$
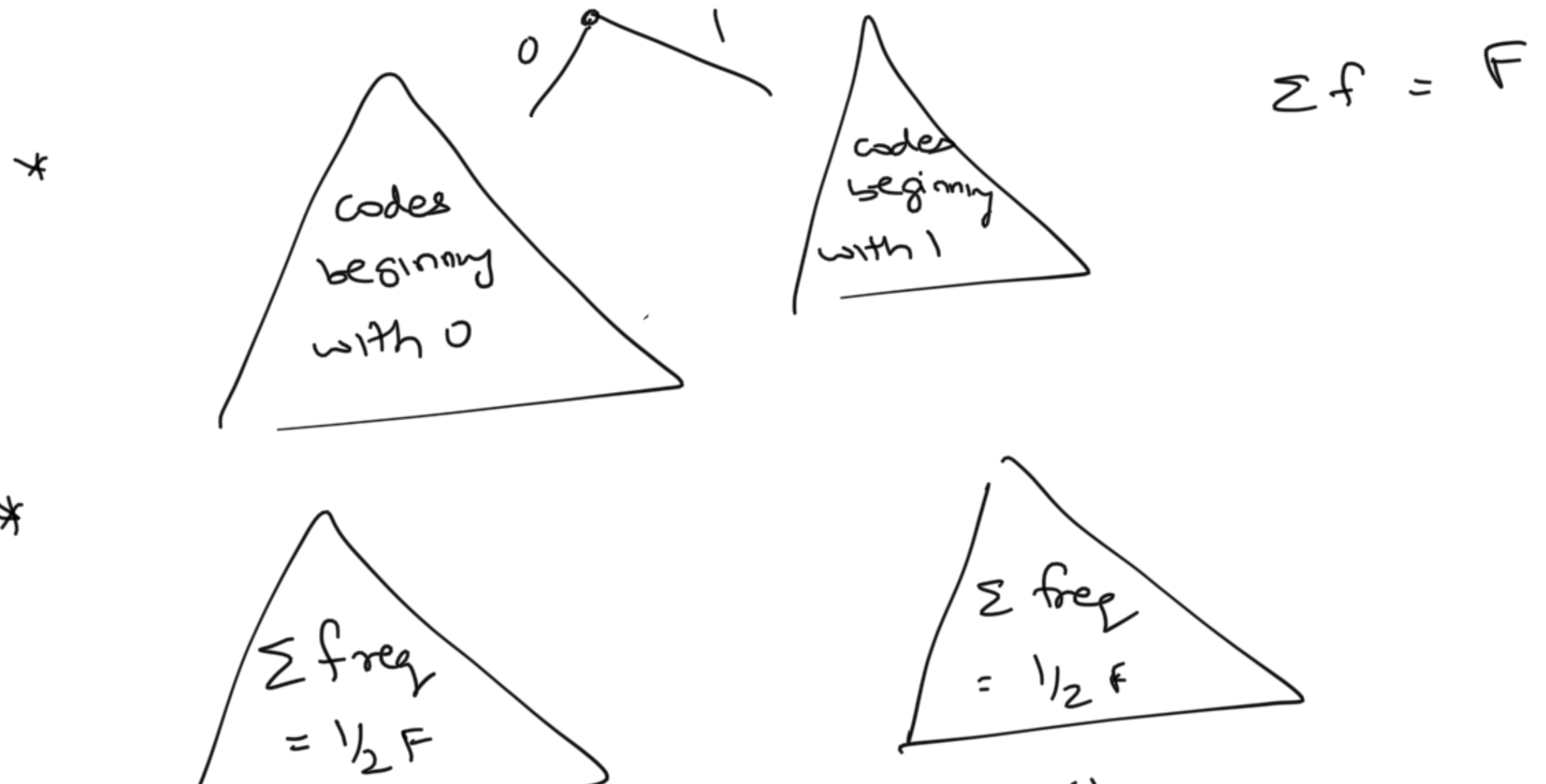$$= O(k^2)$$

→ Can Do BETTER

→ maintain heap instead of array

→ finding min 1, min 2 → $O(\log |A|)$
  + insert new merged freq

→ $O(k \log k)$

\* Greedy: At each stage, choose 2 min freq to be leaves at that stage

Shannon - Fano: Divide and conquer approach
(1950)



$$\Sigma f = F$$

codes beginning with 0

codes beginning with 1

\*

\*

$\Sigma$ freq $= \frac{1}{2} F$

$\Sigma$ freq $= \frac{1}{2} F$

0- subtree
$A_1$

1 - subtree
$A_2$

so $A = A_1 \sqcup A_2$,   $A_1$ total freq $\cong$ $A_2$ total freq

Recursively solve each partition

DOESN'T WORK

Fano's class: Huffman, grad student.
After few years, came up with
greedy sol.